

Itsy Bitsy SpiderNet: Fully Connected Residual Network for Fraud Detection

Sergey Afanasiev
National Research University HSE
Moscow, Russia

Anastasiya Smirnova
National Research University HSE
Moscow, Russia

Diana Kotereva
National Research University HSE
Moscow, Russia

ABSTRACT

With the development of high technology, the scope of fraud is increasing, resulting in annual losses of billions of dollars worldwide. The preventive protection measures become obsolete and vulnerable over time, so effective detective tools are needed. In this paper, we propose a convolutional neural network architecture SpiderNet designed to solve fraud detection problems. We noticed that the principles of pooling and convolutional layers in neural networks are very similar to the way anti-fraud analysts work when conducting investigations. Moreover, the skip-connections used in neural networks make the usage of features of various power in anti-fraud models possible. Our experiments have shown that SpiderNet provides better quality compared to Random Forest and adapted for anti-fraud modeling problems 1D-CNN, 1D-DenseNet, F-DenseNet neural networks. We also propose new approaches for fraud feature engineering called B-tests and W-tests, which generalize the concepts of Benford's Law for fraud anomalies detection. Our results showed that B-tests and W-tests give a significant increase to the quality of our anti-fraud models. The SpiderNet code is available at <https://github.com/aasmirnova24/SpiderNet>

CCS CONCEPTS

• **Security and privacy** → *Usability in security and privacy.*

KEYWORDS

neural networks, fraud detection, CNN, feature engineering

1 INTRODUCTION

The development of high technologies contributes not only to the growth of corporations and world economies but also to the development of fraud, which leads to losses of billions of dollars every year around the world.

In 2018, eight Indian banks incurred \$1.3 billion in losses in a fraud case involving Kingfisher Airlines founder Vijay Mallya¹. In another case, the Agricultural Bank of China faced losses of \$497 million after being defrauded by employees of billionaire Guo Wengui².

Hacker attacks are another global problem. In 2019, the FBI issued an official announcement that global losses from fraudulent Business Email Compromise (BEC) reached \$26 billion during the period from June 2016 to July 2019³.

¹<https://www.theguardian.com/world/2020/apr/20/kingfisher-airlinestycoon-vijay-mallya-loses-appeal-extradition-india>

²<https://www.reuters.com/article/us-china-corruption-tycoonidUSKBN1900DL>

³<https://www.ic3.gov/Media/Y2019/PSA190910>

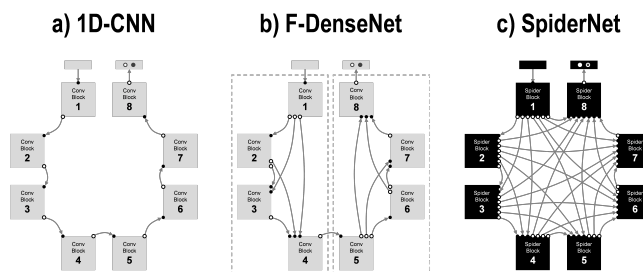


Figure 1: Convolutional neural network architectures designed for fraud detection: a) 1D-CNN, b) F-DenseNet, c) SpiderNet

Another growing threat is social engineering, which has hit Russian bank customers seriously. According to the official data of the Bank of Russia, losses of Russian banks' clients from card fraud reached \$130 million in 2020, which is 10 times higher than similar losses in 2017⁴.

Anti-fraud tools can be roughly divided into directive, preventive, and detective. Directive tools such as instructions and warnings work like a scarecrow and only affect untrained fraudsters. Preventive tools help prevent fraud, but over time, fraudsters adapt and find ways to get around them. Detective tools are essential to detect fraud and minimize losses if fraud has not been prevented. Statistical approaches and machine learning methods are used to develop detective tools. However, there are unresolved problems in this area, such as instability and low generalizing ability of anti-fraud models, as well as high privacy of domain expertise[6].

On the other hand, in recent years, we have seen outstanding advances in deep learning and the successful application of neural networks to practical tasks such as computer vision [21, 22, 34, 62] and natural language processing [19, 49, 51, 57, 58]. This gives us hope that innovative ideas proposed in deep learning will help to remove some of the issues in fraud detection modeling.

In this paper, we propose a convolutional neural network architecture SpiderNet designed to solve fraud detection problems. We noticed that convolutional and pooling layers principles are very similar to the methods of manual processing of information by anti-fraud analysts during investigations. In addition, skip-connections used in convolutional networks [22] make it possible to use features of various power, including fraud scores from external providers.

⁴https://cbr.ru/analytics/ib/fincert/#a_119487

Our proposed technique allows us to increase the generalizing ability of anti-fraud models and is an important advantage of SpiderNet over classical machine learning methods and popular neural network architectures. We show that SpiderNet provides better quality compared to Random Forest and convolutional 1D-CNN and F-DenseNet network architectures adapted for fraud modeling (Figure 1). Moreover, comparing the SpiderNet results with the classic CNN and 1D-DenseNet architectures, we show that the feature locality property is lost while working with tabular data but remains crucial for images. Therefore, the technique of transferring CNN architectures from the CV domain, which is popular in applied tasks, does not give the best result. This confirms the thesis that the application of neural networks to applied tasks requires a deep understanding of the domain area, and it is important to adapt neural network architecture to the specifics of the task.

In addition to the SpiderNet architecture, in this paper, we propose new approaches for developing anti-fraud rules called B-tests and W-tests. Our feature engineering method is based on the idea of identifying manipulation in financial statements using Benford's law [4]. Our proposed approaches can be generalized to any data type, which allows developing B-tests and W-tests for any fraudulent schemes, and not just for identifying accounting manipulations, where Benford's law is applicable [44, 45, 47]. Our results showed that B-tests and W-tests give a significant increase to the quality of our fraud detection models.

To assess the quality of fraud detection models, we use the AUC ROC and Average Precision (AUC PR) metrics, as well as developed by our team business metric PL (Prevented Losses) that allows estimating the funds saved from internal fraud. Our proposed PL metric allows us to solve an important industry issue of assessing the economic efficiency of models developed for industrial use.

We train, evaluate, and compare our SpiderNet with other algorithms on two datasets – private and public. The private dataset contains data on loan applications and internal fraud by POS partners of a large Russian bank, one of the top 50 Russian banks in terms of assets. The public dataset was obtained from an online payment fraud detection competition organized by Ant Financial Services Group⁵ 6.

Testing SpiderNet on two datasets with different types of fraud (internal and transactional) gives us reason to believe that our proposed methods will work well for other types of fraud because SpiderNet is based on the general concepts of fraudulent behavior formulated by Edwin Sutherland and Donald Cressey in criminology [54] and Gary Becker in behavioral economics [3].

The rest of the paper is structured as follows:

- In section 2, we describe a general theoretical model of fraudulent behavior, provide a classification of anti-fraud tools, describe the principles of developing fraud detection models, and outline current problems in this area;
- In section 3, we provide a review of the current state of neural networks for image classification and fraud detection;
- In section 4, we describe the general intuition of the SpiderNet architecture and present the schema of the Spider-Block;

- In Section 5, we describe methods for the automated development of B-tests and W-tests anti-fraud rules;
- In Section 6, we describe our experiment: provide characteristics of datasets, feature engineering methods, data preprocessing methods, tricks for training models, hyperparameters of algorithms, and metrics for models' quality assessment;
- In section 7, we demonstrate the results of the experiments;
- And in section 8, we summarize, draw general conclusions and outline open questions for future research.

2 FRAUD DETECTION

The foundations of modern criminology and the theory of white-collar fraud were laid nearly 100 years ago by Edwin Sutherland and his student Donald Cressey [54], who are considered some of the most influential criminologists of the 20th century. Sutherland and Cressey proposed to consider the crime not from the position of criminal law, as was customary, but from the position of sociology, applying basic sociological concepts in criminology.

Four decades later, Nobel laureate Gary Becker, using the principle of economics imperialism, described the economic model of crime [3], according to which crime can be viewed as an activity that some people choose rationally, comparing the expected benefits and expected costs:

$$(1 - \pi) * U(W_C) - \pi * S > U(W_L) \quad (1)$$

where π – the probability of being caught (assessed by the criminal, i.e. subjectively);

$U(\cdot)$ – individual utility function;

S – penalties incurred in the event of capture (for example, a fine or criminal punishment);

W_C -- proceeds of crime;

W_L -- income from legal activities.

The left side of inequality (1) characterizes crime-related elements; the right side is the utility from legal earnings. Logically, when inequality is satisfied, the individual, other things being equal, will prefer to break the law.

To minimize losses from fraudulent activities, companies need to develop and implement anti-fraud tools that affect the components of inequality (1). In particular, the total losses from fraud can be reduced by increasing the probability π and decreasing the W_C component by increasing fraudsters' costs for bypassing anti-fraud protection.

Anti-fraud actions can be divided into 3 types⁷:

- *Directive* – instructions, regulations, training materials, contracts, etc.;
- *Preventive* – tools that are aimed at preventing fraudulent activities (locks, safes, passwords, etc.);
- *Detective* – tools used to detect fraud (anti-fraud rules and models, investigation techniques, etc.).

Preventive tools are the most effective, but fraud is highly adaptive, which means that fraudsters constantly come up with new ways how to get around the company's security. That is why the

⁵<https://dc.cloud.alipay.com/index#/topic/data?id=4>

⁶<https://www.kaggle.com/gmhost/atec-anti-fraud/version/2>

⁷The Information Security Standards (CISSP) use more granular typing: preventative, deterrent, detective, corrective, recovery, compensation, directive, administrative, logical/technical, physical.

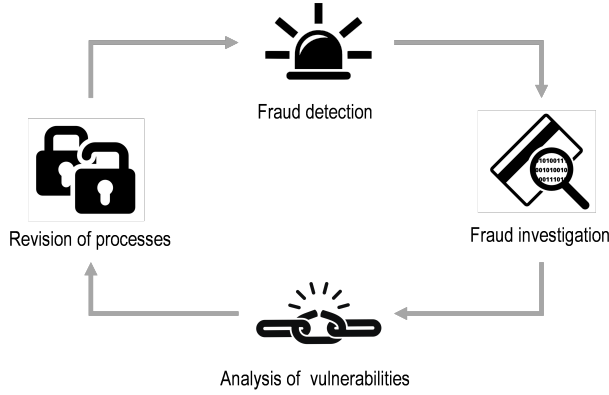


Figure 2: Arrangement of anti-fraud processes in the bank according to the principles of the scientific method.

fight against fraud is compared to "confrontation of armor and projectiles". To face this constant arms race, the company's anti-fraud sector requires an integrated, systems approach. One of the main principles of the system approach is the arrangement of anti-fraud processes in the form of a cyclic scheme according to the principles of the scientific method.

The scheme in Figure 2 shows how anti-fraud technologies are developed and improved. One of the main ideas of the scheme is that fraud and counteraction to fraud (anti-fraud) constantly influence each other [59]. It means that, according to the scheme, the company should regularly review its anti-fraud processes, modifying and improving them. Therefore, detective tools that allow detecting new fraudulent schemes and vulnerabilities in processes are an important part of the anti-fraud system in a company.

The detective tools are based on strong rules developed by experts and anti-fraud models built on these strong rules. Strong rules development consists of three key stages (Figure 3):

- (1) At the first stage, simple features and rules are developed by experts;
- (2) At the second stage, complex rules are combined using arithmetic and logical operations over simple features and rules;
- (3) At the third stage, strong rules with high predictive power for detecting fraud are selected from complex rules.

The general principle of developing strong rules is similar to the method of collecting evidence in forensic science when all kinds of evidence are first collected, and then, based on the combinations of the collected evidence, a crime is proved. A similar principle is at the heart of convolutional neural networks, in which sequential convolution operations form complex features, and pooling operations filter out noise, i.e. select strong features that are good at predicting the target variable. This is why convolutional networks can be considered a powerful and intuitive algorithm for developing fraud detection models.

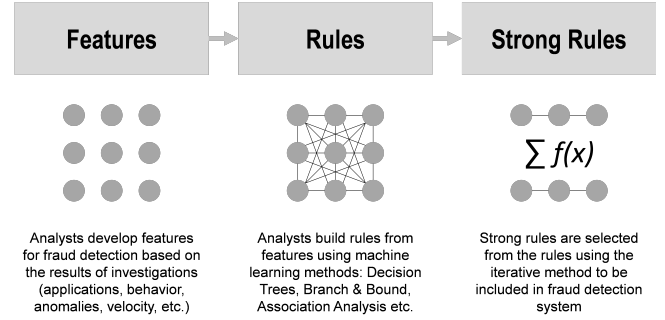


Figure 3: Scheme of developing strong rules for fraud detection.

3 RELATED WORK

3.1 CNN Architectures

Over the past decade, convolutional neural networks (CNNs) have made a breakthrough in computer vision problems solution [50]. The basic principles of CNN's work were borrowed from the works of Hubel and Wiesel, who studied the visual cortex in the 1950s and 1960s [30, 31]. The first CNN architecture was proposed by Yann LeCun in the late 1980s [66], and in the late 1990s, LeCun's research group developed the LeNet-5 architecture [37], which consisted of convolutional and pooling layers that perform the function of implicit regularization of the neural network.

A breakthrough in computer vision came in 2012 when the deep convolutional neural network AlexNet won the ILSVRC computer vision competition in image classification problems [34]. The authors noticed that increasing CNN depth has an implicit regularization effect, which later became mainstream and was exploited in such architectures as VGG [53], NiN [42], Inception, GoogleNet [55], ResNet [22], and others.

In 2014, the Google team showed excellent results at ILSVRC, proposing the GoogLeNet architecture [55], a feature of which was Inception blocks with bottlenecks and 1x1 convolutions, which allowed increasing the number of convolutional channels, i.e. the width of the neural network. This technique also became popular and found its application in the architectures WRN [69], Xception [10], ResNeXt [65], MobileNets [27], NASNet [71, 72], etc.

In 2015, the Microsoft Research team proposed another breakthrough architecture – ResNet [22] that used skip connections, which create an ensemble effect. This technique proved to be very effective and was later used in WRN [69], Xception [10], ResNeXt [65], FractalNet [36], DenseNet [28], etc.

In 2019, the Google Research team published an EfficientNet approach [56] to scale deep CNNs. Authors noticed that the key characteristics of convolutional architectures, such as depth, width, and resolution, depending on each other, therefore, for scaling, it is necessary to select the optimal combination of these parameters. Later, the EfficientNet approach allowed the participants to take the

top places in the Deepfake Detection Challenge on Kaggle, which was organized by Facebook in early 2020.

In addition to the developed architectures, effective tricks have been proposed for training CNN:

- Regularization method Weight Decay (L2 regularization) [23, 32, 35];
- The Dropout technique [26], which allows to accidentally disconnect links in fully connected layers;
- Stochastic Depth technique [29] used in ResNet networks to random disable blocks;
- Augmentation methods: Cutout [14], Mixup [70], CutMix [68], Valid/Diverse Noise, Flipping, etc. [64];
- Learning-rate reduction strategies that improve the convergence of gradient descent methods [16, 40, 43, 48].

These and other proposed ideas allowed to achieve high results in Computer Vision.

3.2 Fraud Prediction Models

Statistical methods and machine learning have been used for fraud detection tools development for many years. In 2002, Bolton and Hand [6] highlighted the key problems of using statistical methods in fraud detection modeling. Here are some of these issues that remain relevant today:

- (1) The scope of fraud is increasing with the development of high technology;
- (2) Fraudsters bypass preventive technologies over time, so detective tools are needed;
- (3) On the other hand, detective algorithms also degrade over time, so they need to be regularly updated;
- (4) Databases and anti-fraud methods are closed from the scientific community. This makes it difficult to research and develop this area;
- (5) To develop anti-fraud models, unsupervised (search for anomalies) and supervised (search for known fraudulent patterns) algorithms are used. Unsupervised models make a lot of mistakes, because anomalies may be caused by operational errors, marketing promotions, etc. Supervised models are trained on historical data, so they are poor at catching new fraudulent schemes;
- (6) There is a global problem of class imbalance - there are much fewer fraudulent observations than non-fraudulent ones. This leads to a high false positive rate, which is why many false-positives are sent for investigation and the use of models in anti-fraud processes becomes expensive.

In the discussion of this paper, Provost and Breiman noted other important issues:

- (7) Models are customized for a specific fraudulent scheme, which makes them difficult to scale to other types of fraud;
- (8) Big data is needed to develop effective anti-fraud models, so this remains the lot of large companies;
- (9) Machine learning algorithms do not solve the anti-fraud problem, expert knowledge and understanding of fraudulent schemes are needed.

The development of anti-fraud models is most often solved as a binary classification problem (fraud/non-fraud), where classical

algorithms are used, such as SVM, Logistic Regression, Random Forest, and others [5].

With the onset of the deep learning boom, neural networks began to be used in fraud detection modeling and almost completely replaced classical machine learning methods in papers [33].

Wiese and Omlin [63] proposed using a recurrent LSTM network to detect fraudulent credit card transactions. The authors suggested that since recurrent networks were designed to process sequences, they should be good at detecting fraudulent patterns in card transaction sequences. Authors were able to demonstrate the advantage of the LSTM over the SVM algorithm, but the LSTM network failed to beat the simple fully connected neural network FFNN due to the insufficient number of fraudulent transactions in the dataset.

Fu et al. [17] solved a similar problem in detecting card fraud by training the architecture of the LeNet-5 convolutional neural network developed by Yann LeCun for image processing. To train LeNet-5 on the card fraud detection problem, the authors presented transactions in the form of rectangular matrices of features, which were fed to the CNN input as two-dimensional pictures. The results obtained showed the superiority of LeNet-5 over the classical algorithms SVM, ANN, and Random Forest.

Heryadi and Warnars [25] continued to develop these ideas and tried to combine CNN with a recurrent LSTM network. The authors' hypothesis was the following: CNNs, due to convolutions, should detect short-term fraudulent patterns, LSTM network, due to long short-term memory, should work well with long sequences of fraudulent transactions. Authors have developed a hybrid CNN-LSTM architecture, but experiments have shown that simple CNN detects fraud better than CNN-LSTM. Authors made an important conclusion from their results: long-term fraudulent schemes that cannot be detected for a long time are extremely rare, in contrast to short-term fast fraudulent schemes. This is confirmed by Becker's economic model of crime (eq. 1).

Attempts to apply the popular neural network architectures CNN and RNN for fraud detection tasks continued in subsequent research efforts.

Li et al. [39] used the DenseNet architecture to detect electricity theft in China.

Chen and Liu [8] refined the DenseNet architecture by adding an Inception module to the beginning of the network and additional skip connections between the Inception layer and Dense blocks. Their new CNN architectures named LI and DI have improved the results of standard CNN and DenseNet for the task of transaction fraud detection.

Cheng et al. [9] proposed a Spatio-temporal neural network STAN based on attention. The STAN architecture includes an Attention module and a simple CNN. For the task of detecting transaction fraud, STAN showed better quality compared to CNN, LSTM, etc.

Li and Liu [41] proposed to use a special loss function for transaction fraud detection, which solved the problem of intraclass variability. Their loss function FCL (full center loss), constructed as a combination of DCL (distance center loss) and ACL (angle center loss), worked like batch normalization.

For the tasks of organized fraud on Internet sites (fake reviews, bots, spam, etc.) detection, graph neural networks are gaining popularity today, allowing feature extraction for interconnected objects [15, 60, 61, 67].

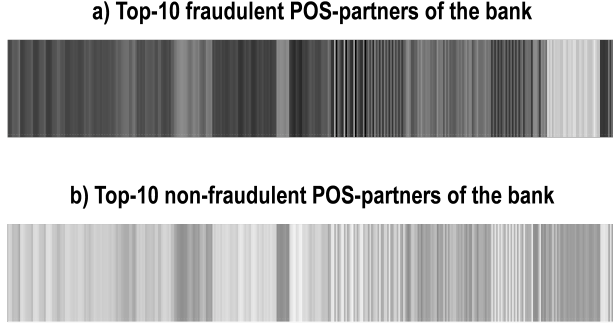


Figure 4: Examples of one-dimensional heatmaps of features for internal fraud detection: a) average feature values for the top-10 fraudulent POS-partners of the bank; b) average feature values for the top-10 non-fraudulent POS-partners of the bank.

4 SPIDERNET

4.1 Problem Formulation

One of the main problems of using neural networks in fraud detection tasks is that many of the proposed architectures were migrated from other domains (mainly from popular CV and NLP) without a deep understanding of why these architectures should work on fraud detection tasks.

When designing a neural network architecture for fraud detection, our intuition is that anti-fraud rules developed by experts are a kind of digital evidence (by analogy with pixels in images, which are digital features processed by convolutions in CNN). At the same time, anti-fraud rules have different power like forensic evidence. Moreover, anti-fraud rules can work in conjunction with each other, strengthening the evidence base. This intuition tells us that CNN’s convolution and pooling operations are the most appropriate tools for combining anti-fraud rules and selecting strong combinations of them.

In computer vision problems, it was shown that a sequence of several convolutional layers allows one to create hierarchical feature maps using the locality property in images [38]. On the other hand, the locality property is lost in tabular data, since the features can be arranged in a different order. Nevertheless, studies of CNN architectures have shown that convolutions learn well not only medium-frequency (local) features but also low-frequency ones (texture, background, color, etc.) [18, 24]. This allows us to assume that CNN architectures should perform well on tabular data, where 1D convolutions can be applied to the feature vector. In classification problems, such feature vectors will differ for observations of different classes in color and texture (Figure 4).

The mathematical intuition of how convolutions work can be explained through CNN regularization by drop connections between neurons. Moreover, convolutions drop connections in an orderly manner, as opposed to the random dropout used in fully connected networks (Figure 5). In addition, the receptive field in convolutions

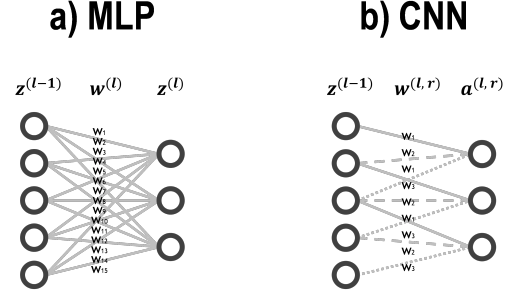


Figure 5: The dropout principle between neurons in CNN: a) connections between layers in a fully connected neural network (MLP); b) connections between layers in a CNN.

allows capturing a small number of features, due to which many small models are trained in CNN on these subsets, which are then summed up, achieving the ensemble effect.

On the other hand, if the combination of anti-fraud rules obtained on the hidden layers has a strong predictive ability, then we want to use this combination without additional processing, forwarding it to the output layer of the neural network using skip connections.

This intuition well represents the best practices of anti-fraud investigations and can be implemented using a fully connected residual network, which we call SpiderNet (Figure 1).

4.2 Spider Block

SpiderNet consists of blocks that are connected using skip connections. Thus, each block receives features from all previous blocks, and these features are processed using convolutional and pooling layers and are forwarded to all subsequent blocks. Cause we want to have only the strongest features at the output, the blocks that are farthest from the network output contain several pooling layers, filtering out weak and medium features. Conversely, the closer the block is to the network output, the fewer pooling layers it contains since the features that have reached this block have already passed several convolutional and pooling layers. The general architecture of SpiderNet’s block is shown in Figure 6.

SpiderNet is a convolutional feedforward network with skip connections between blocks. Formally, the k th Spider-block is defined by the recursive formula:

$$y_k = \mathcal{F}_k (y_{k-1} \oplus \dots \oplus y_1) = \mathcal{F}_k \left(\sum_{i=1}^{k-1} \oplus y_i \right) \quad (2)$$

where y_k is a vector of $n - k$ outputs for the k th block;

$\mathcal{F}_k(\cdot)$ is the k th block operator, which combines the functions dropout, convolution, batch normalization, ReLU, and Max-pooling;

\oplus is the concatenation operator for incoming vectors;

$\sum \oplus$ is a short notation for vector concatenation;

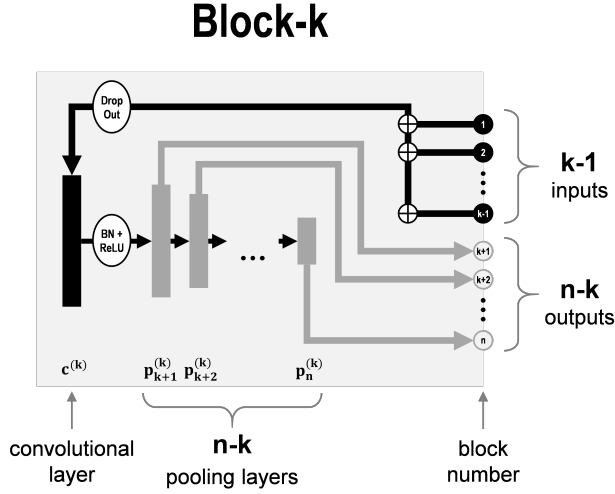


Figure 6: Scheme of the k th Spider-block with convolutional layer and $n-k$ pooling layers (n is the total number of Spider-blocks).

y_i is the output vector of the i th block, which is fed to the input of the k th block ($1 \leq i < k$).

In the last n th Spider-block, after convolutional layer and Batch Norm+ReLU and MaxPooling layer, there is also global average pooling which is applied to reduce the dimension of the channels. After these operations, the vector y_n is fed to two fully connected layers with dropout and SoftMax output for binary classification.

Expression (2) demonstrates the mathematical intuition of the SpiderNet architecture. The output vector of the k th block is obtained by transforming the concatenated outputs from the previous blocks, which are smaller neural networks, which receive concatenated vectors of the previous blocks as input, etc. Thus, SpiderNet works as an ensemble of neural networks, which allows improving the convergence and generalization of the entire neural network.

5 B-TESTS AND W-TESTS

The peculiarity of fraud modeling is that most of the fraud rules are developed manually by experts. This because fraudulent schemes are diverse, especially internal fraud schemes that are designed for specific vulnerabilities of the organization.

We offer B-tests and W-tests approaches that automate the manual process and generalize feature engineering for various types of internal fraud.

5.1 B-tests

B-tests are based on Benford's law, discovered at the end of the 19th century by Simon Newcomb [46], who noticed that the values of some numerical data start with one more often than two, two more often than three, and so on. Later, Frank Benford [4] empirically confirmed this law for various social and physical phenomena, and the first-digit law was named after him.

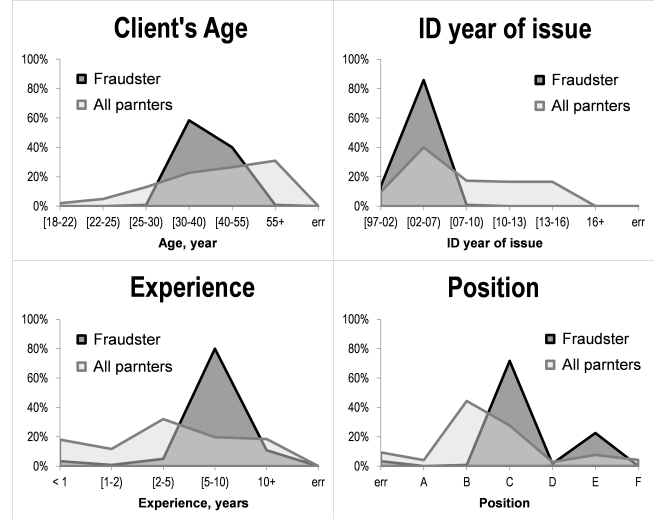


Figure 7: Examples of B-tests for internal fraud detection.

In the 90s of the twentieth century, Mark Nigrini used Benford's law to audit financial statements and managed to identify managerial embezzlement of \$2 million in the office of the Arizona State Treasurer [47].

Summarizing these ideas for all types of data (numerical and categorical), we propose a B-tests technique, which consists of comparing the distributions of data characterizing an object with the general population of all objects (for example, the activities of employees or partners of the company).

The general idea is that internal fraud is rare, i.e. it does not greatly affect the distribution of the general population, while the data on the activities of fraudsters are very different from the average (Figure 7).

B-tests can be customized according to various characteristics of the sample, such as the partner sphere of business, the region of the partner's location, the number of quantiles in distributions, the divergence cut-off threshold, etc. Thus, the setting of B-tests can be reduced to an optimization task of searching for anomalies [1].

B-tests can be calculated using various metrics reflecting the divergence between two distributions, for example, Chi-squared test, K-S test, Anderson-Darling test, etc.

For our B tests, we calculate the area difference for two discrete distributions using the following formula:

$$S = \frac{1}{2} \sum_{i=1}^n |a_i - b_i| \quad (3)$$

where a_i and b_i are the compared distributions;

n is the number of quantiles in the distribution.

The number of quantiles in the distribution n and the threshold for S depend on the number of objects in the samples and are tunable hyperparameters.

5.2 W-tests

W-tests solve the same problem as B-tests using the Wasserstein metric, which is defined as:

$$W_p(\mu, \nu) = \inf_{\gamma \in \Gamma(\mu, \nu)} E_{(x,y) \sim \gamma} [\|x - y\|] \quad (4)$$

where $E[Z]$ denotes the expected value of a random variable Z and the infimum is taken over all joint distributions of the random variables X and Y with marginals μ and ν respectively.

Wasserstein metric solves the problem of the insufficient number of objects in samples, which is typical for B-tests. However, the Wasserstein metric can be calculated only for numeric data types, so W-tests cannot completely replace B-tests.

6 EXPERIMENT SETUP

6.1 Datasets

We trained and compared SpiderNet with other algorithms on two datasets: private data and public data. The general characteristics of datasets are presented in Table 1.

Private dataset. Our private dataset contains data on the credit activity of POS partners of a large Russian bank, one of the top 50 Russian banks in terms of assets. Before applying the models, the bank used a rule-based approach. When the rules were triggered, the top-worst POS partners were sent to the bank’s security for investigation. Since the security staff is limited, the number of investigations is also limited and amounts to 40 investigations per week. Thus, the main goal of applying the model approach was to reduce the time it takes to identify fraudulent partners with an unchanged investigation budget.

The bank has a procedure for blocking unprofitable POS partners based on risk indicators, which are calculated approximately 90 days after the first fraudulent loan is issued. Therefore, to form the target variable, the loss for the POS partner was calculated within 90 days from the date of the fraud rules calculation. Unprofitable POS partners were marked as fraudulent, profitable ones - as non-fraudulent. This approach to constructing the target variable makes it possible to detect problem partners.

Since fraud rules could have been triggered during the loss assessment period, the evaluation of the model approach will show the uplift to the existing rule-based process, and not the full economic effect of the model.

The records in the private dataset are presented as a vector of features for each POS partner as of the weekly slice date. Therefore, a single POS partner can be included in the sample several times with different slice dates. The depth of calculation of the features ranges from 7 to 60 days ago from the date of the slice. A flag, indicating an internal fraud event, was set to each record as a target variable. The overall task comes down to predicting the internal fraud of the POS partner based on the data of its historical activities.

Public dataset. The public dataset comes from an online payment fraud detection competition hosted by Ant Financial Services Group. The dataset contains the values of features for payment transactions and a binary target variable that reflects whether the payment is fraudulent or not.

6.2 Feature Selection

To form the final samples, we did data preprocessing. Firstly, we checked the features for validity and removed features with a low

Table 1: Characteristics of private and public datasets

| Attribute | Private Data | Public Data |
|------------------------------|-----------------|------------------------------|
| Source | Russian Bank | Ant Financial Services Group |
| Type of data | POS credits | Payments |
| Type of fraud | Internal | Transaction |
| Period time | 03.2014-10.2019 | 09.2017-11.2017 |
| All observations, # | 1 880 499 | 990 006 |
| Fraud observations, # | 5 327 | 12 122 |
| Fraud rate, % | 0.28 | 1.22 |
| All features, # | 509 | 297 |
| Selected features, # | 163 | 128 |

fill rate. Secondly, we selected features using a cross-correlation matrix. Of a pair of highly correlated features, we left the one that had a stronger correlation with the target variable.

As a result of data preprocessing, 163 features remained in the private dataset, and 128 features remained in the public dataset.

6.3 Model Architectures and Tricks

To assess the generalization ability and demonstrate the benefits of SpiderNet, we compared the results of several machine learning algorithms:

- (1) *Random Forest* is a baseline model that has been chosen as a strong algorithm and industry standard for modeling on tabular data. We used 5-fold cross-validation and the Optuna library to tune the Random Forest hyperparameters [2].
- (2) *1D-CNN* is a one-dimensional convolutional network with classical alternation of convolutional and pooling layers and two fully connected layers and a SoftMax layer. We trained CNN with 3, 6, and 8 convolutional layers;
- (3) *1D-DenseNet* is a classic DenseNet architecture [28] with an implementation for one-dimensional vectors. We trained two-block architectures with 3 and 4 convolutions in each block;
- (4) *F-DenseNet* is a DenseNet architecture adapted for fraud prediction with two fully connected convolutional blocks containing convolutional and pooling layers. We trained architectures with 3 and 4 convolutional layers in each block;
- (5) *SpiderNet* is our fully connected residual convolutional network with convolutional-pooling blocks. We trained 6 and 8 block architectures (Figure 8).

We used weight decay (L2-regularization) in BatchNorm and fully connected layers, and dropout on fully connected layers as regularizers for all neural networks. In the 1D-CNN, F-DenseNet, and SpiderNet architectures we applied the BatchNorm and ReLU transformations after each convolutional layer. In the F-DenseNet and SpiderNet blocks, as an additional regularization for the skip connections, we used dropout after concatenating the input vectors for each block (Figure 4). We also used the fraud-rate leveling technique in batches to solve the problem of the lack of fraud samples in batches in case of over-class imbalance.

To train and assess the quality of the models, we performed a stratified split of private and public datasets into a train (80%),

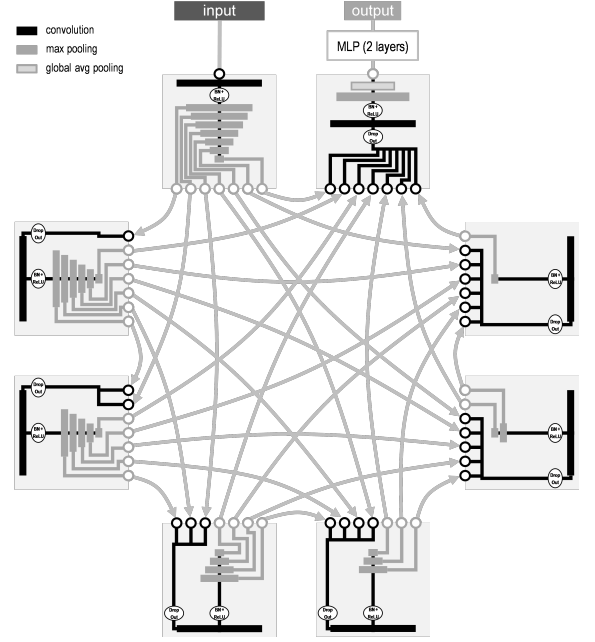
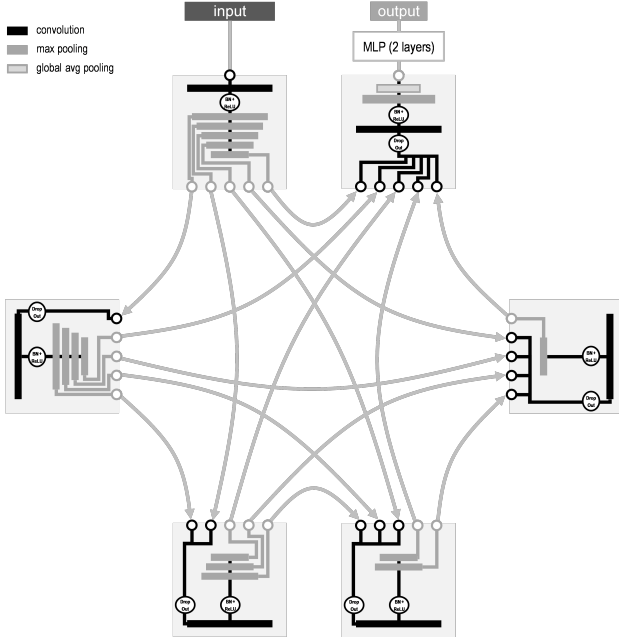


Figure 8: SpiderNet-6 (left) and SpiderNet-8 (right) neural network architectures for fraud detection tasks.

validation (10%), and test (10%) samples. We tuned the networks to the validation sample using Grid-Search and early stopping.

6.4 Performance Measures

We used AUC ROC and AUC PR metrics to evaluate the quality of models. The AUC ROC metric and the linearly related Gini coefficient are industry standards in banking modeling. However, as shown by Saito and Rehmsmeier [52], the AUC ROC accepts unreasonably high values and becomes uninformative on over unbalanced samples, in contrast to the AUC PR, which adequately estimates the quality of models on unbalanced samples. Therefore, we tuned the hyperparameters of the models using the AUC PR.

To evaluate the quality of fraud detection models on the private dataset, we have developed the special metric PL (Prevented Loss), which shows how much loss from internal fraud the model prevents. To calculate the PL, we estimate prevented loss for each partner:

$$PL^{(i)} = P_{(T_l - T_a)} \cdot \frac{DR - DR_0}{1 - DR_0} \quad (5)$$

where T_l is the whole considered period of loss;

T_a is the period on which the model works;

$(T_l - T_a)$ is the period after the model is triggered (for our sample, it is 90 days – the empirical period for which the bank’s Security detects fraud without using the model);

DR is the Default-Rate for loans issued by the partner for the period $(T_l - T_a)$;

DR_0 is a "zero target" for Default-Rate in which the loan portfolio has zero profit;

$P_{(T_l - T_a)}$ is the partner’s loan portfolio for the period $(T_l - T_a)$.

Total PL is calculated based on the company’s total investigative resources, i.e. on the assumption that security costs do not increase.

Total Prevented Loss shows the net profit from the model due to the earlier detection of fraud than it had been before the model was applied:

$$PL = \sum_{i=1}^k PL^{(i)} \cdot b_i \quad (6)$$

where $PL^{(i)}$ is prevented loss for the i th fraud partner;

k is the number of the first k partners with the highest model probability of fraud (for our bank $k = 40$);

b_i is a binary variable showing the event for the i th partner: 1 – fraud, 0 – no fraud.

For the public dataset, the economic metric was not developed due to the lack of the necessary financial data in the dataset.

7 EXPERIMENT

7.1 Model Optimization and Tuning

We tuned models by AUC PR and selected the best hyperparameters for the models.

Random Forest:

Private data: class_weight=0.0167, max_depth=7, n_estimators=129;

Public data: class_weight=0.0744, max_depth=7, n_estimators=90.

CNN:

Private data: layers=8, l2_batch=0.0001, kernel_size=3, dropout=0.25, weight_decay=0.0002, filters=10, learning_rate=0.003, hidden=100;

Public data: layers=8, l2_batch=0.0001, kernel_size=5, dropout=0.25, weight_decay=0.0002, filters=10, learning_rate=0.003, hidden=30.

DenseNet:

Private data: block_sizes = [4, 4], initial_filters=5, initial_stride=1,

Table 2: Quality of models for internal (private dataset) and transactional (public dataset) fraud detection: the best results are highlighted in bold; 95% confidence intervals are shown in parentheses.

| # | Model | Private data (test sample) | | Public data (test sample) | |
|----|---------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| | | AUC PR | AUC ROC | AUC PR | AUC ROC |
| 1 | Random Forest | 0.0650 (± 0.001116) | 0.9371 (± 0.009253) | 0.4881 (± 0.003114) | 0.9709 (± 0.003572) |
| 2 | CNN-3 | 0.0527 (± 0.001012) | 0.9339 (± 0.012978) | 0.4462 (± 0.003096) | 0.9670 (± 0.004674) |
| 3 | CNN-6 | 0.0644 (± 0.001111) | 0.9385 (± 0.011432) | 0.4908 (± 0.003114) | 0.9711 (± 0.004780) |
| 4 | CNN-8 | 0.0708 (± 0.001161) | 0.9288 (± 0.009605) | 0.5099 (± 0.003114) | 0.9718 (± 0.004511) |
| 5 | DenseNet-6 [3; 3] | 0.0646 (± 0.001113) | 0.9315 (± 0.009091) | 0.4757 (± 0.003111) | 0.9669 (± 0.004935) |
| 6 | DenseNet-8 [4; 4] | 0.0691 (± 0.001148) | 0.9310 (± 0.010545) | 0.4854 (± 0.003113) | 0.9686 (± 0.004661) |
| 7 | F-DenseNet-6 [3; 3] | 0.0732 (± 0.001179) | 0.9263 (± 0.014509) | 0.5092 (± 0.003114) | 0.9708 (± 0.005082) |
| 8 | F-DenseNet-8 [4; 4] | 0.0575 (± 0.001054) | 0.9186 (± 0.015820) | 0.4968 (± 0.003114) | 0.9704 (± 0.004780) |
| 9 | SpiderNet-6 | 0.0948 (± 0.001326) | 0.9484 (± 0.008004) | 0.5375 (± 0.003106) | 0.9721 (± 0.004763) |
| 10 | SpiderNet-8 | 0.0680 (± 0.001139) | 0.9277 (± 0.009588) | 0.5160 (± 0.003113) | 0.9684 (± 0.004744) |

k=5, conv_kernel_width=3, bottleneck_size=2, theta=0.5, transition_pool_stride=1, initial_conv_width=5, initial_pool_width=2, initial_pool_stride=2;

Public data: block_sizes = [4, 4], initial_filters=5, initial_stride=1, k=10, conv_kernel_width=5, bottleneck_size=2, theta=0.5, transition_pool_stride=1, initial_conv_width=5, initial_pool_width=2, initial_pool_stride=2.

F-DenseNet:

Private data: blocks = 2, convolutions=3+3, weight_decay=0.0002, l2_batch=0.0002, kernel_size=7, filters=15, hidden=60, dropout=0.25, learning_rate=0.003;

Public data: blocks = 2, convolutions=3+3, weight_decay=0.0002, l2_batch=0.0002, kernel_size=5, filters=15, hidden=100, dropout=0.25, learning_rate=0.003.

SpiderNet:

Private data: blocks=6, l2_batch=0.0001, kernel_size=3, filters=10, hidden = 100, weight_decay=0, dropout=0.25, learn_rate=0.005, dropout_block_k=0.001 * k^4 (where k={4, 5} is a block number);

Public data: blocks=6, l2_batch=0.0002, kernel_size=7, filters=15, hidden=30, weight_decay=0.0002, dropout=0.25, learn_rate=0.005, dropout_block_k=0.001 * k^4 (where k={4, 5} is a block number).

For all neural networks, we used the Adam optimizer.

7.2 Results

The results of the trained models showed that the SpiderNet-6 architecture had the best quality on both datasets for both AUC PR and AUC ROC (Table 2). Confidence intervals were calculated using the asymptotic method for AUC PR [7] AUC ROC [13]. We see clear dynamics that increase in the number of skip-connections results in increase in quality for fraud detection models: CNN-8: 0 skip connections \rightarrow F-DenseNet-6: 2 skip connections \rightarrow SpiderNet-6: 10 skip connections. This confirms our hypothesis that strong features can pass well from the different layers of neural network immediately to its output due to skip connections. On the other hand, we see that the best model of the classic DenseNet-8 performs worse than the best CNN-8, which has no skip connections. This indirectly confirms the hypothesis that, for tabular data, bottlenecks between blocks don't allow strong features to go directly to the

Table 3: Quality of models for internal fraud detection.

| # | Model | Private data (test sample) | |
|----|---------------------|----------------------------|--------------------|
| | | Fraud, # | PL |
| | Random classifier | 48 | \$325 604 |
| 1 | Random Forest | 208 | \$2 079 527 |
| 2 | CNN-3 | 312 | \$2 235 707 |
| 3 | CNN-6 | 280 | \$2 753 821 |
| 4 | CNN-8 | 280 | \$2 337 297 |
| 5 | DenseNet-6 [3; 3] | 240 | \$2 324 181 |
| 6 | DenseNet-8 [4; 4] | 288 | \$2 433 914 |
| 7 | F-DenseNet-6 [3; 3] | 240 | \$2 297 848 |
| 8 | F-DenseNet-8 [4; 4] | 272 | \$2 402 470 |
| 9 | SpiderNet-6 | 304 | \$2 570 014 |
| 10 | SpiderNet-8 | 264 | \$2 379 977 |
| | Perfect classifier | 888 | \$4 659 439 |

Table 4: Sign tests for two pairs of results (the PL and Fraud metrics are normalized according to the perfect classifier).

| | CNN-3 | Spider Net-6 | CNN-6 | Spider Net-6 |
|----------------------|---------------|---------------|---------------|---------------|
| Private data: | | | | |
| AUC PR | 0.0527 | 0.0948 | 0.0644 | 0.0948 |
| AUC ROC | 0.9339 | 0.9484 | 0.9385 | 0.9484 |
| PL (recall) | 0.4798 | 0.5516 | 0.5910 | 0.5516 |
| Fraud (recall) | 0.3514 | 0.3423 | 0.3153 | 0.3423 |
| Public data: | | | | |
| AUC PR | 0.4462 | 0.5375 | 0.4908 | 0.5375 |
| AUC ROC | 0.9670 | 0.9721 | 0.9711 | 0.9721 |
| p-value | 0.015625 | | 0.015625 | |

output layers of the network. The F-DenseNet architecture, which is adapted for fraud detection tasks, does not contain a bottleneck between blocks, so the quality of the best F-DenseNet-6 surpasses the quality of CNN-8 and DenseNet-8. These results demonstrate the importance of developing a neural network architecture for the specifics of the task.

We also confirm Saito’s and Rehmsmeier’s results [52] that the AUC ROC metric is not informative for problems with unbalanced data and shows unreasonably high values. For our datasets, the AUC PR metric is preferred. Although the results for the private dataset according to the AUC PR metric are low, we consider them satisfactory, since they don’t take into account the full effect of the models, but demonstrate an uplift when switching from the rule-based approach to model one. In addition, the evaluation of these models according to the PL metric showed a significant positive economic effect with an unchanged budget for investigations (Table 3). Since the PL metric is not normalized, Table 3 also shows the efficiency for the random and perfect classifiers on the private dataset. The results demonstrate that the SpiderNet has an almost 8-fold increase in financial efficiency compared to a random classifier and has less than a 2-fold decrease to the perfect model (as if the bank had identified all fraudulent POS partners before committing crimes – so-called "Minority Report").

We also see that for the private dataset, SpiderNet-6 ranked second according to PL metric, losing to the simpler CNN-6. This may be because the hyperparameters in the Grid-Search were selected following the integral metric AUC PR, while the PL metric is a threshold and covers a small number of POS-partners who have been checked by the bank’s security (40 partners per week out of 6786 on average).

On the other hand, as we can see from the results of Table 3, SpiderNet-6 outperforms CNN-6 by the number of identified fraudulent POS partners (out of all limited number of those verified by the security) but loses to CNN3 by this indicator. The number of fraudulent partners detected is also an important metric of the model quality. Hand et al. [20] suggest that in mass fraud, it is better to tune the model to maximize the number of detected fraud cases than to maximize prevented losses. This approach can have a greater influence on fraudsters, since, as we know from Becker’s concept, the number of identified fraudulent cases directly proportionally affects the probability of being caught (eq. 1). The results obtained show that on various datasets and various quality metrics, SpiderNet-6 loses to CNN-3 and CNN-6 only in one metric. At the same time, in 5 other cases, SpiderNet-6 wins over CNN-3 and CNN-6. To assess the statistical significance of the obtained SpiderNet-6 advantage, we tested the null hypothesis and alternative one-sided hypothesis using a sign test [11]:

- (1) For SpiderNet-6 and CNN-3

$$H_0^{(1)} : \text{models are of the same quality}$$

$$P(\text{Quality}_{\text{SpiderNet-6}} > \text{Quality}_{\text{CNN-3}}) = 0.5$$

$$H_1^{(1)} : \text{SpiderNet-6 has higher quality}$$

$$P(\text{Quality}_{\text{SpiderNet-6}} > \text{Quality}_{\text{CNN-3}}) > 0.5$$

- (2) For SpiderNet-6 and CNN-6

$$H_0^{(2)} : \text{models are of the same quality}$$

$$P(\text{Quality}_{\text{SpiderNet-6}} > \text{Quality}_{\text{CNN-6}}) = 0.5$$

$$H_1^{(2)} : \text{SpiderNet-6 has higher quality}$$

$$P(\text{Quality}_{\text{SpiderNet-6}} > \text{Quality}_{\text{CNN-6}}) > 0.5$$

The obtained results of p-values are presented in Table 4 and we see that in both pairwise comparisons the null hypothesis is

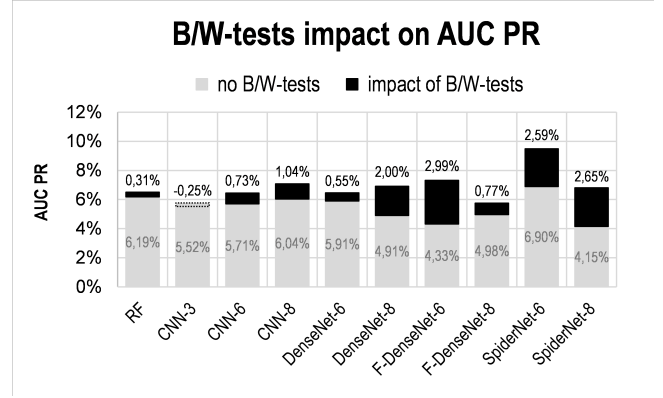


Figure 9: Influence of B-tests and W-tests on model quality (negative impact means a decrease in the model quality when adding B/W-tests)

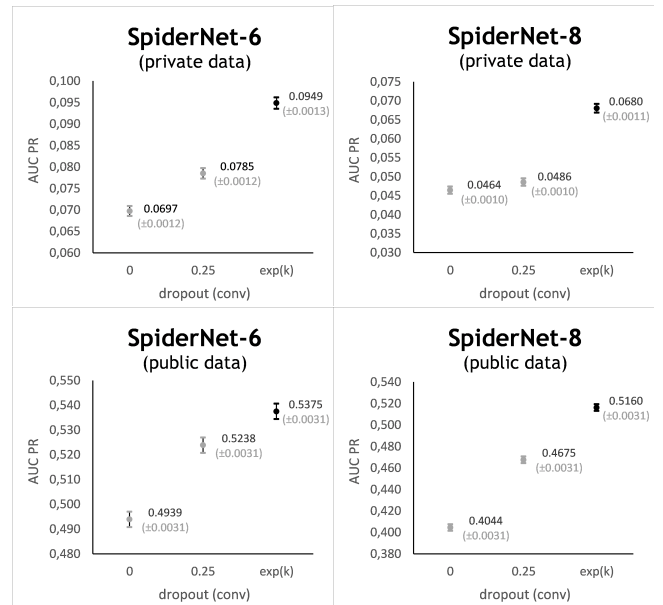


Figure 10: SpiderNet quality (AUC PR), depending on the dropout techniques in convolutional layers: 1) the value of zero corresponds to SpiderNet quality without a dropout; 2) the value of 0.25 corresponds to constant dropout=0.25 used in all Spider-blocks (vanilla technique); 3) the value of exp(k) corresponds to an exponential increase in the dropout value as the Spider-block number increases (our technique).

rejected in favor of the alternative, i.e. with a 95% significance level SpiderNet-6 more often shows better quality than CNN-3 or CNN-6.

We also tested the effectiveness of B-tests and W-tests by training the models on two private samples:

- (1) Full sample containing 24 B-tests, 15 W-tests, and 124 expert rules;
- (2) Truncated sample containing 124 expert rules.

Our results showed the high efficiency of B-tests and W-tests by AUC PR (Figure 9). Adding B-tests and W-tests gave us the AUC PR increase of 19.4%. At the same time, we see that B/W-tests have a greater impact on models with skip connections. This may be because B/W-tests are strong rules and for a better impact they must be sent immediately to the network output.

We also see that SpiderNet-6 trained on expert rules without B-tests shows the best quality by AUC PR metric, which once again proves the stability of the results obtained.

In addition, the results of experiments demonstrated that our exponential dropout technique in convolutional layers provides significant gains in quality over a constant or zero dropout (Figure 10). This trick allows adding more regularization in the last blocks, which get more information from the previous blocks.

All program codes and detailed results are available at: <https://github.com/aasmirnova24/SpiderNet>

8 CONCLUSIONS

In this paper, we investigated deep learning methods for fraud prediction and proposed a novel architecture of neural networks for fraud detection problems. Taking inspiration from the skip connection concepts of Resnet [22], FractalNet [36], Adanet [12], and DenseNet [28] convolutional networks for computer vision, we developed SpiderNet architecture for fraud detection, guided by anti-fraud expert knowledge. Using convolutional layers, our SpiderNet creates hierarchical combinations of anti-fraud rules, and a fully connected structure of skip connections between blocks allows strong rules to be forwarded immediately to the network output. Also, our network can select strong rules early on through the use of a multi-layered structure of pooling layers in Spider-blocks. Moreover, the down-dropout technique we use between Spider-blocks is an additional regularizer against the excessive complexity of the network.

In our opinion, SpiderNet should work well for heterogeneous input data, when there are clear leaders among the rules supplied to the network input and they must be forwarded to the output without additional transformation (for example, scores of other models or strong rules).

It can also be noted that our results confirm the hypothesis of ResNet authors, which they formulated in their paper [22]: "The residual learning principle is generic, and we expect that it is applicable in other vision and non-vision problems".

In this paper, we also proposed new methods for developing anti-fraud rules – B-tests and W-tests, which significantly affect the quality of the models. We hope that B/W-tests will contribute to the solution of the problem "The Blind Men and the Elephant", which was formulated by Foster Provost [6].

We also hope that the metric of quality Prevented Losses developed by us allows us to solve an important issue for the industry: how to evaluate the economic efficiency of anti-fraud models developed for industrial use.

We understand that our SpiderNet does not solve all the anti-fraud modeling problems identified by Bolton, Hand, Provost, and Breiman [6]. In particular, we still use expert rules designed for specific fraud types to train models. Our B-tests and W-tests partially solve this problem, but there are other strong methods for fraud

feature engineering, such as graph methods [15, 60, 61, 67], entropy changing methods [17], and variance anomaly detection methods (V-tests, similar to B-tests). We plan to work on these topics in our future research.

An important component of SpiderNet is skip-connection, which helps to forward strong features directly to the output layers of the network, partially solving the problem of locality in convolutions, when the order of features in the input vector is important, and their rearrangement leads to a change in the quality of the model. However, the current implementation of SpiderNet does not completely solve the locality problem. Our future work will focus on this problem.

We also assume that SpiderNet might work well not only for fraud detection but also for other types of modeling tasks that use tabular data. Testing this hypothesis is also a topic for our future research.

ACKNOWLEDGMENTS

We are grateful to Dmitry Efimov for his valuable advice.

REFERENCES

- [1] Sergey Afanasiev and Anastasiya Smirnova. 2018. Predictive fraud analytics: B-tests. *Journal of Operational Risk* 13 (Dec. 2018). <https://doi.org/10.21314/JOP.2018.213>
- [2] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. *Optuna: A Next-generation Hyperparameter Optimization Framework*. <http://arxiv.org/abs/1907.10902>
- [3] Gary S. Becker. 1968. Crime and punishment: an economic approach. *Journal of Political Economy* 76 (2) (March 1968), 169–217. <https://doi.org/10.1086/259394>
- [4] Frank Benford. 1938. The law of anomalous numbers. *Proc. Am. Philos. Soc.* 78 (4) (March 1938). <https://doi.org/10.2307/984802>
- [5] Siddhartha Bhattacharyya, Sanjeev Jha, Kurian Tharakunnel, and J. Christopher Westland. 2011. *Data mining for credit card fraud: A comparative study*. <https://doi.org/10.1016/j.dss.2010.08.008>
- [6] Richard J. Bolton and David Hand. 2002. Statistical Fraud Detection: A Review. *Statist. Sci.* 17 (3) (March 2002), 235–255. <https://doi.org/10.1214/ss/1042727940>
- [7] Kendrick Boyd, Kevin Hasegawa Eng, and C. David Page. 2013. Area Under the Precision-Recall Curve: Point Estimates and Confidence Intervals. *ECML PKDD* (2013). https://doi.org/10.1007/978-3-642-40994-3_29
- [8] Zhuo Chen and Guanjun Liu. 2019. *DenseNet+Inception and Its Application for Electronic Transaction Fraud Detection*. <https://doi.org/10.1109/HPCC/SmartCity/DSS.2019.00357>
- [9] Dawei Cheng, Sheng Xiang, Chencheng Shang, Yiyi Zhang, Fangzhou Yang, and Liqing Zhang. 2020. *Spatio-Temporal Attention-Based Neural Network for Credit Card Fraud Detection*. <https://doi.org/10.1609/aaai.v34i01.5371>
- [10] Francois Chollet. 2017. *Xception: Deep Learning with Depthwise Separable Convolutions*. <https://arxiv.org/abs/1610.02357>
- [11] W.J. Conover. 1999. *Practical Nonparametric Statistics*. Wiley, New York.
- [12] Corinna Cortes, Xavi Gonzalvo, Vitaly Kuznetsov, Mehryar Mohri, and Scott Yang. 2017. *AdaNet: Adaptive Structural Learning of Artificial Neural Networks*. <https://arxiv.org/abs/1607.01097>
- [13] Corinna Cortes and Mehryar Mohri. 2004. Confidence Intervals for the Area under the ROC Curve. <https://papers.nips.cc/paper/2004/file/a7789ef8d599b8df86bbe632b2994d-Paper.pdf>
- [14] Terrance DeVries and Graham W. Taylor. 2017. *Improved regularization of convolutional neural networks with cutout*. <https://arxiv.org/abs/1708.04552>
- [15] Yingdong Dou, Zhiwei Liu, Li Sun, Yutong Deng, Hao Peng, and Philip S. Yu. 2020. *Enhancing Graph Neural Network-based Fraud Detectors against Camouflaged Fraudsters*. <https://doi.org/10.1145/3340531.3411903>
- [16] Stanislav Fort and Stanislaw Jastrzebski. 2019. *Large Scale Structure of Neural Network Loss Landscapes*. <https://arxiv.org/abs/1906.04724>
- [17] Kang Fu, Dawei Cheng, Yi Tu, and Liqing Zhang. 2016. Credit Card Fraud Detection Using Convolutional Neural Networks. *Lecture Notes in Computer Science* 9949 (2016). https://doi.org/10.1007/978-3-319-46675-0_53
- [18] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and harnessing adversarial examples. <https://arxiv.org/abs/1412.6572>
- [19] Alex Graves, Abdel rahman Mohamed, and Geoffrey Hinton. 2013. Speech recognition with deep recurrent neural networks. *IEEE international conference on acoustics, speech and signal processing* 9949 (2013), 6645–6649. <https://arxiv.org/abs/1303.5778>

- [20] David Hand, Chris Whitrow, Niall M. Adams, Piotr Juszczak, and Dave Weston. 2008. Performance criteria for plastic card fraud detection tools. *Journal of the Operational Research Society* (2008). <https://doi.org/10.1057/palgrave.jors.2602418>
- [21] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision* (2017), 2961–2969. <https://arxiv.org/abs/1703.06870>
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep Residual Learning for Image Recognition. *Microsoft Research* (2015). <https://arxiv.org/pdf/1512.03385v1.pdf>
- [23] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. 2018. Mask r-cnn. *Bag of Tricks for Image Classification with Convolutional Neural Networks* (2018). <https://arxiv.org/abs/1812.01187>
- [24] Katherine L. Hermann, Ting Chen, and Simon Kornblith. 2020. The Origins and Prevalence of Texture Bias in Convolutional Neural Networks. <https://arxiv.org/abs/1911.09071>
- [25] Yaya Heryadi and Harco Leslie Hendric Spits Warnars. 2017. Learning temporal representation of transaction amount for fraudulent transaction recognition using CNN, Stacked LSTM, and CNN-LSTMs. *IEEE International Conference on Cybernetics and Computational Intelligence (CyberneticsCom)* (2017). <https://doi.org/10.1109/CYBERNETICSCOM.2017.8311689>
- [26] Geoffrey E. Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R. Salakhutdinov. 2012. *Improving neural networks by preventing co-adaptation of feature detectors*. <https://arxiv.org/abs/1207.0580>
- [27] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. 2017. *MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications*. <https://arxiv.org/abs/1704.04861>
- [28] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. 2018. Densely Connected Convolutional Networks. (2018). <https://doi.org/10.1109/CVPR.2017.243>
- [29] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Weinberger. 2016. *Deep Networks with Stochastic Depth*. <https://arxiv.org/abs/1603.09382>
- [30] D.H. Hubel and T.N. Wiesel. 1959. Receptive Fields of Single Neurons in the Cat's Striate Cortex. *Journal of Physiology* 148 (1959), 574–591. <https://doi.org/10.1113/jphysiol.1959.sp006308>
- [31] D.H. Hubel and T.N. Wiesel. 1962. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *Physiol* 160 (1962), 106–154. <https://doi.org/10.1113/jphysiol.1962.sp006837>
- [32] Xianyan Jia, Shutao Song, Wei He, Yangzhaohao Wang, Haidong Rong, Feihu Zhou, Liqiang Xie, Zhenyu Guo, Yuanzhou Yang, Liwei Yu, Tiegang Chen, Guangxiao Hu, Shaohuai Shi, and Xiaowen Chu. 2018. *Highly scalable deep learning training system with mixed-precision: Training imagenet in four minutes*. <https://arxiv.org/abs/1807.11205>
- [33] Kanika and Jimmy Singla. 2020. *A Survey of Deep Learning based Online Transactions Fraud Detection Systems*. <https://doi.org/10.1109/ICIEM48762.2020.9160200>
- [34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
- [35] Anders Krogh and John A. Hertz. 1992. A simple weight decay can improve generalization. , 950–957 pages. <http://papers.nips.cc/paper/563-a-simple-weight-decay-can-improve-generalization.pdf>
- [36] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. 2016. *FractalNet: Ultra-Deep Neural Networks without Residuals*. <https://arxiv.org/abs/1605.07648>
- [37] Yann Lecun, Leon Bottou, Y. Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86 (11) (1998), 2278–2324. <https://doi.org/10.1109/5.726791>
- [38] Honglak Lee, Roger Grosse, Rajesh Ranganath, and Andrew Y. Ng. 2009. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. *ICML* (2009). <https://doi.org/10.1145/1553374.1553453>
- [39] Bo Li, Kele Xu, Xiaoyan Cui, Yiheng Wang, Xinbo Ai, and Yanbo Wang. 2018. Multi-Scale DenseNet-Based Electricity Theft Detection. https://www.researchgate.net/publication/344779530_Multi-scale_DenseNet-Based_Electricity_Theft_Detection
- [40] Yuanzhi Li, Colin Wei, and Tengyu Ma. 2020. Towards Explaining the Regularization Effect of Initial Large Learning Rate in Training Neural Networks. (2020). <https://arxiv.org/abs/1907.04595>
- [41] Zhenchuan Li, Guanjin Liu, and Changjun Jiang. 2020. Deep Representation Learning with Full Center Loss for Credit Card Fraud Detection. *IEEE Transactions on Computational Social Systems* (2020). <https://doi.org/10.1109/TCSS.2020.2970805>
- [42] Min Lin, Qiang Chen, and Shuicheng Yan. 2014. Network In Network. (2014). <https://arxiv.org/abs/1312.4400>
- [43] Ilya Loshchilov and Frank Hutter. 2017. SGDR: Stochastic Gradient Descent with Warm Restarts. *ICLR* (2017). <https://arxiv.org/abs/1608.03983>
- [44] Fletcher Lu and J. Efrim Boritz. 2005. Detecting Fraud in Health Insurance Data: Learning to Model Incomplete Benford's Law Distributions. In *16th European Conference on Machine Learning* (2005), 633–640. https://doi.org/10.1007/11564096_63
- [45] Fletcher Lu, J. Efrim Boritz, and Dominic Covey. 2006. Adaptive Fraud Detection Using Benford's Law. *Advances in Artificial Intelligence* (June 2006), 347–358. https://doi.org/10.1007/11766247_30
- [46] Simon Newcomb. 1881. Note on the frequency of use of the different digits in natural numbers. *American Journal of Mathematics* 4 (1) (1881), 39–40.
- [47] Mark J. Nigrini. 1999. I've Got Your Number: How a mathematical phenomenon can help CPAs uncover fraud and other irregularities. *Journal of Accountancy* (1999).
- [48] Brendan O'Donoghue and Emmanuel Candes. 2017. Adaptive Restart for Accelerated Gradient Schemes. <https://doi.org/10.1007/s10208-013-9150-3>
- [49] Romain Paulus, Caiming Xiong, and Richard Socher. 2017. A deep reinforced model for abstractive summarization. <https://arxiv.org/abs/1705.04304>
- [50] Maithra Raghu and Eric Schmidt. 2020. *A Survey of Deep Learning for Scientific Discovery*. <https://arxiv.org/abs/2003.11755>
- [51] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. <https://arxiv.org/abs/1606.05250>
- [52] Takaya Saito and Marc Rehmsmeier. 2015. The Precision-Recall Plot Is More Informative than the ROC Plot When Evaluating Binary Classifiers on Imbalanced Datasets. *PLoS ONE* 10(3) (2015). <https://doi.org/10.1371/journal.pone.0118432>
- [53] Karen Simonyan and Andrew Zisserman. 2015. Very deep convolutional networks for large-scale image recognition. (2015). <https://arxiv.org/pdf/1409.1556.pdf>
- [54] Edwin H Sutherland and Donald R Cressey. 1992. *Principles of Criminology*. AltaMira Pres, Lanham, Md.
- [55] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2014. *Going deeper with convolutions*. <https://arxiv.org/pdf/1409.4842v1.pdf>
- [56] Mingxing Tan and Quoc V. Le. 2019. *EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks*. <https://arxiv.org/abs/1905.11946>
- [57] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. 2016. *Wavenet: A generative model for raw audio*. <https://arxiv.org/abs/1609.03499>
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. <https://arxiv.org/abs/1706.03762>
- [59] Walter L. Wallace. 1971. The Logic of Science and Sociology. (1971). <https://doi.org/10.4324/97813151329763>
- [60] Daixin Wang, Jianbin Lin, Peng Cui, Quanhuai Jia, Zhen Wang, Yanming Fang, Quan Yu, Jun Zhou, Shuang Yang, and Yuan Qi. 2019. A Semi-supervised Graph Attentive Network for Financial Fraud Detection. (2019). <https://doi.org/10.1109/ICDM.2019.00070>
- [61] Haobo Wang, Zhao Li, Jiaming Huang, Pengrui Hui, Weiwei Liu, Tianlei Hu, and Gang Chen. 2020. Collaboration Based Multi-Label Propagation for Fraud Detection. (2020). <https://doi.org/10.24963/ijcai.2020/343>
- [62] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. 2018. Non-local neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2018), 7794–7803. <https://arxiv.org/abs/1711.07971>
- [63] Benard Wiese and Christian Omlin. 2007. Credit Card Transactions, Fraud Detection, and Machine Learning: Modelling Time with LSTM Recurrent Neural Networks. *Innovations in Neural Information Paradigms and Applications* 247 (2007), 231–268. <https://doi.org/10.1007/978-3-642-04003-0>
- [64] Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-Thang Luong, and Quoc V. Le. 2019. Unsupervised data augmentation. (2019). <https://arxiv.org/abs/1904.12848>
- [65] Saining Xie, Ross Girshick, Piotr Dollar, Zhuowen Tu, and Kaiming He. 2017. Aggregated Residual Transformations for Deep Neural Networks. (2017). <https://doi.org/10.1109/CVPR.2017.634>
- [66] B. Boser Y. LeCun, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. 1989. Backpropagation applied to handwritten zip code recognition. *Neural computation* (1989), 2278–2324. <https://doi.org/10.1162/neco.1989.1.4.541>
- [67] Shuhan Yuan, Xintao Wu, Jun Li, and Aidong Lu. 2019. Spectrum-based deep neural networks for fraud detection. (2019). <https://arxiv.org/abs/1706.00891>
- [68] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. 2019. CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features. (2019). <https://arxiv.org/abs/1905.04899>
- [69] Sergey Zagoruyko and Nikos Komodakis. 2016. Wide Residual Networks. *British Machine Vision Conference* (2016). <https://doi.org/10.5244/C.30.87>
- [70] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. 2017. mixup: Beyond Empirical Risk Minimization. (2017). <https://arxiv.org/abs/1710.09412>
- [71] Barret Zoph and Quoc V. Le. 2017. Neural Architecture Search with Reinforcement Learning. <https://arxiv.org/abs/1611.01578>
- [72] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. 2018. Learning Transferable Architectures for Scalable Image Recognition. <https://arxiv.org/abs/1707.07012>