



POLITECNICO
MILANO 1863

**SCUOLA DI INGEGNERIA INDUSTRIALE
E DELL'INFORMAZIONE**

Credit card fraud detection

NUMERICAL ANALYSIS FOR MACHINE LEARNING

COMPUTER SCIENCE AND ENGINEERING (T2I - ARTIFICIAL INTELLIGENCE)

First Name: Davide

Family Name: Gesualdi

Student Number: 101761 - Personal Code: 10885255

Professor:
Edie MIGLIO

Course credits:
10

Academic year:
2023-2024

Abstract: The project aims to implement a method for credit card fraud detection, building upon the research presented in "Efficient credit card fraud detection using evolutionary hybrid feature selection and random weight networks" [1]. The proliferation of digital transactions has concurrently increased the vulnerability of credit card systems to fraudulent activities. Effective fraud detection is crucial to safeguard financial assets. The paper introduces an efficient method for credit card fraud detection that relies on Competitive Swarm Optimization (CSO) and Random Weight Network (RWN). Additionally, the system includes an automated hybrid feature selection capability to identify the most pertinent features during the detection process. The experimental outcomes validate that this system can attain outstanding results in G-Mean, RUC, and Recall values.

1. Introduction

Credit card fraud detection has become a critical challenge in modern financial transactions, driven by the rising prevalence of online transactions and digital payment methods. The unauthorized use of credit cards for fraudulent activities presents considerable financial risks for individuals, businesses and financial institutions. Machine learning algorithms have revolutionized the field of fraud detection, providing advanced capabilities for automated data analysis and pattern recognition [2] [3]. These algorithms can learn from historical transaction data and identify anomalies that may signal fraudulent behavior. The problem of credit card fraud detection is particularly complex from a machine learning perspective. Notably, the dataset is inherently imbalanced, with legitimate transactions vastly outnumbering fraudulent ones. Moreover, the concept of fraud evolves over time as consumer habits and fraudulent tactics undergo changes, resulting in concept drift within the data. The selection of relevant features from transaction data is a cornerstone of constructing effective fraud detection models. Not all attributes contribute equally to the discrimination between legitimate and fraudulent transactions [4]. Feature selection is the process of identifying and retaining those attributes that are most influential in detecting fraudulent patterns, while discarding irrelevant or redundant information. Methods like filter and wrapper approaches are commonly employed to carry out feature selection, allowing models

to achieve enhanced performance, characterized by improved accuracy, reduced false positive rates and better generalization across diverse fraud scenarios^[5]. The proposed efficient credit card fraud detection model, named HybridIG-CSO, utilizes an automatic feature selection mechanism that identifies significant features through the combined use of Information Gain (IG) and Competitive Swarm Optimization (CSO) techniques, with a Random Weight Network (RWN) serving as the foundational classifier. The objective of this project is to implement the HybridIG-CSO model and subsequently compare its performance against several fundamental classifiers, including Naïve Bayes (NB), Random Forest (RF) and Support Vector Machine (SVM).

2. Dataset Analysis

The models were evaluated using four distinct datasets, as outlined in the table below, where "Abb." denotes the assigned dataset code and "#PS" signifies the positive samples in each dataset:

Dataset	Abb.	#Samples	#Features	#PS	Link
Loan Prediction	D ₁	614	11	192 (31%)	https://github.com/Paliking/ML_examples/blob/master/LoanPrediction/train_u6lujuX_CVtuZ9i.csv
Creditcardcsvpresent	D ₂	3075	10	448 (14%)	https://github.com/gksj7/creditcardcsvpresent/blob/main/creditcardcsvpresent.csv
Default ofCreditCardClients	D ₃	30000	23	6636 (22%)	https://archive.ics.uci.edu/dataset/350/default+of+credit+card+clients
European cardholders	D ₄	284807	30	482 (0.17%)	https://kaggle.com/mlg-ulb/creditcardfraud

Table 1: The dataset characteristics.

2.1. Skewness

Skewness is a statistical measure that quantifies the asymmetry of the distribution of values in a dataset. It indicates whether the data points are skewed to the left (negative skew) or the right (positive skew) relative to the mean.

- Positive Skewness (Right Skew): the majority of data points are concentrated on the left side of the distribution, with a few extreme values on the right side.
Mean > Median > Mode
- Negative Skewness (Left Skew): most of the data points are concentrated on the right side of the distribution, while a few extreme values on the left side.
Mean < Median < Mode
- Zero Skewness (Symmetrical Distribution): perfectly symmetrical distribution, where the mean, median, and mode are equal. In such a distribution, the data points are evenly distributed around the central point.

Skewness was computed only for numerical attributes containing no missing values.

```
1 skewed_feats = dataset[numerical_cols].apply(lambda x: skew(x.dropna()))
2 skewness = pd.DataFrame({"Variable":skewed_feats.index, "Skewness":skewed_feats.values})
```

Variable	Skewness
ApplicantIncome	6.523526
CoapplicantIncome	7.473217
Credit_History	-1.877351
LoanAmount	2.670763
Loan_Amount_Term	-2.356504

Table 2: Skewness dataset D₁.

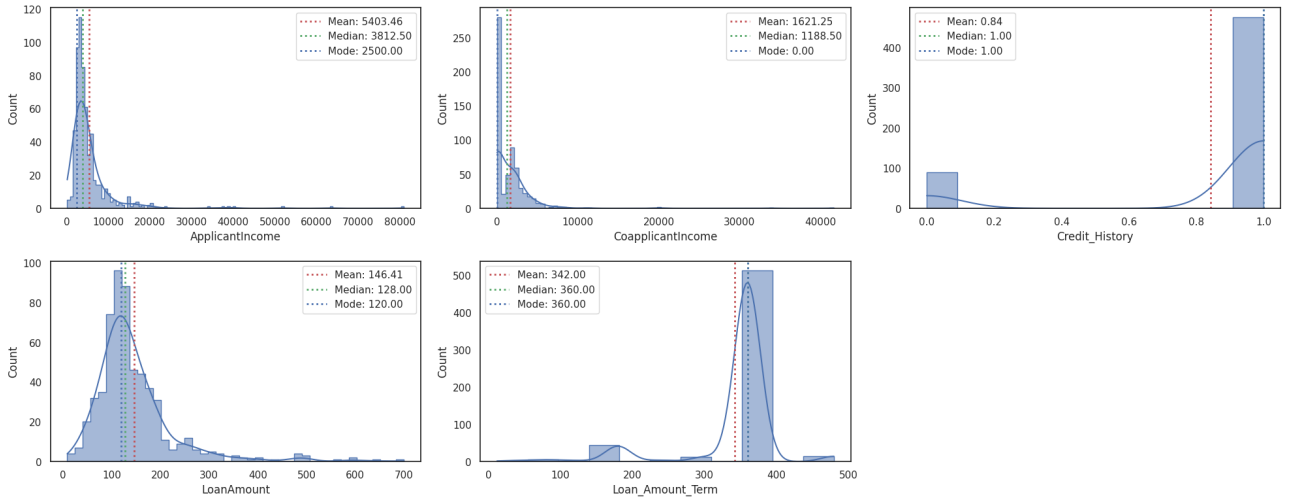


Figure 1: Histograms to show distributions of dataset D₁ with a kernel density estimate to smooth the distribution and vertical lines denoting the mean (red), median (green), mode (blue).

2.2. Correlation Analysis

The computation of pairwise attribute correlations using the Pearson correlation coefficient (excluding missing values) revealed that, across all four datasets, the attributes are predominantly weakly linear correlated, either positively or negatively. Only a small number of attributes exhibit moderate positive linear correlations. The exception is the third dataset, which exhibits a higher number of strong linear correlations compared to the other datasets. Nevertheless, these strong correlations are still relatively modest in comparison to the overall prevalence of very weak correlations.

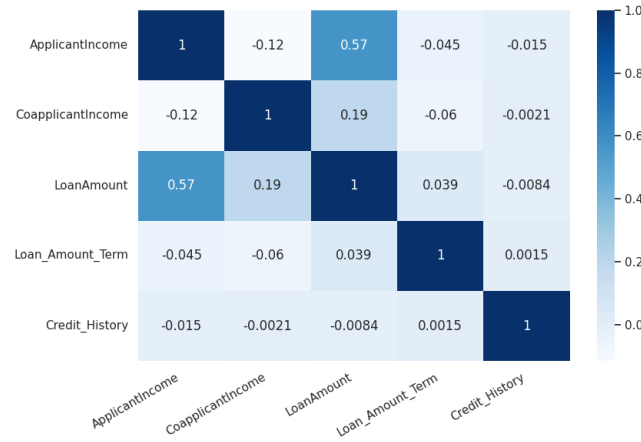


Figure 2: Correlation Analysis of dataset D₁.

2.3. Dataset splitting

The dataset was shuffled and split into training, validation and test sets, with proportions of 70%, 15% and 15%, respectively. During the testing phase, model performance was assessed both on a separate test set and using a 10-fold cross-validation approach.

2.4. Handling missing values

Missing values were exclusively present in the first dataset. Categorical features underwent label encoding prior to handling missing data.

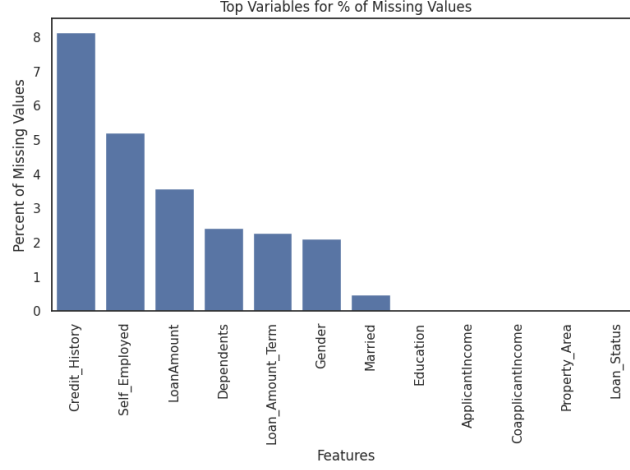


Figure 3: Top Variables by percentage of Missing Values in the dataset D_1 .

To address missing values, a K-Nearest Neighbors (KNN) imputation approach was employed. For numerical features, KNN mean imputation was implemented using the `KNNImputer` class from the `scikit-learn` library, with the number of neighbors set to 3. Consequently, missing values in numerical features for each sample were imputed with the mean value of its K nearest neighbors in the training dataset.

```

1 X_train_imputed = X_train.dropna(subset=categorical_cols)
2 imputer = KNNImputer(n_neighbors=n_neighbors).fit(X_train_imputed)
3 X_train_imputed = pd.DataFrame(
4     imputer.transform(X_train_imputed),
5     columns=X_train.columns,
6     index=X_train_imputed.index)

```

Categorical features were imputed using KNN majority vote imputation through the `scikit-learn` `KNeighborsClassifier` class.

```

1 def majority_vote_imputation(df, df_imputed, classifier, categorical_cols, n_neighbors):
2     for col in categorical_cols:
3         to_impute = df[df[col].isna() & df.drop(columns=col).notna().all(axis=1)].copy()
4         for i in range(len(to_impute)):
5             sample = to_impute.drop(columns=col).iloc[i]
6             to_impute[col].iloc[i] = classifier[col].predict([sample])[0]
7         df_imputed = pd.concat([df_imputed, to_impute])
8     return df_imputed

```

Samples containing missing values in multiple categorical features were removed from the dataset, resulting in a reduction of dataset size of only 13 samples (2%).

2.5. Data normalization

Box plots and violin plots generated from the original datasets revealed a clear need for data normalization, due to the pervasive presence of outliers and the disparate scales across variables.

After evaluating the `StandardScaler`, the `QuantileTransformer` was selected for data normalization.

`StandardScaler` is sensitive to outliers, as it relies on the mean and standard deviation for scaling, which can be significantly influenced by extreme values.

In contrast, `QuantileTransformer` mitigates the impact of outliers by using quantiles information. `QuantileTransformer` applies a non-linear transformation such that the probability density function of each feature will be mapped to a uniform or Gaussian distribution. In this case, all the data, including outliers, will be mapped to a uniform distribution within the range $[0, 1]$, making outliers indistinguishable from inliers.

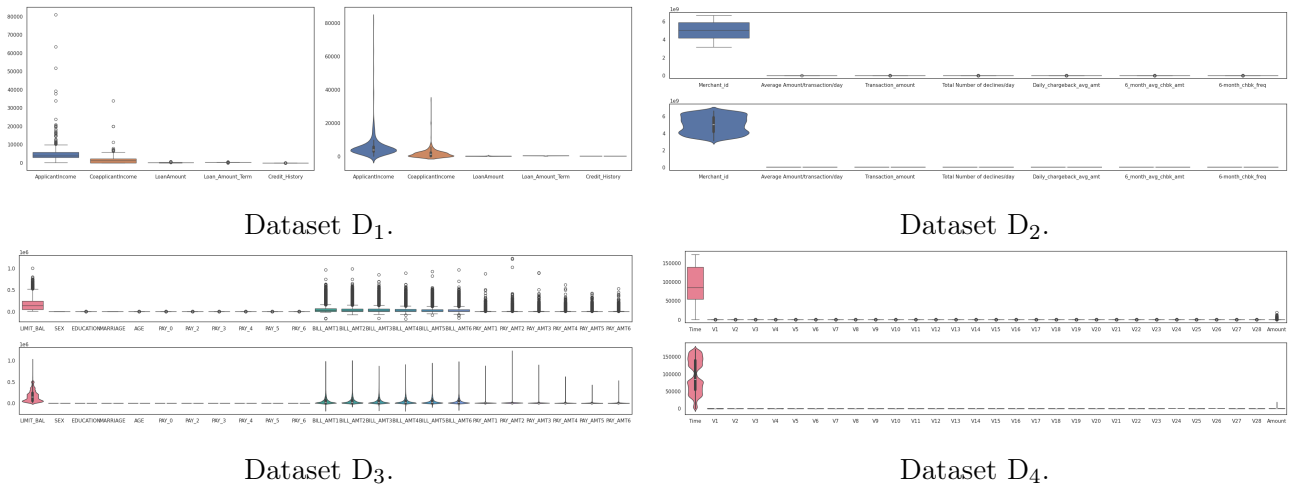


Figure 4: Box plots and violin plots of the original data in the datasets.

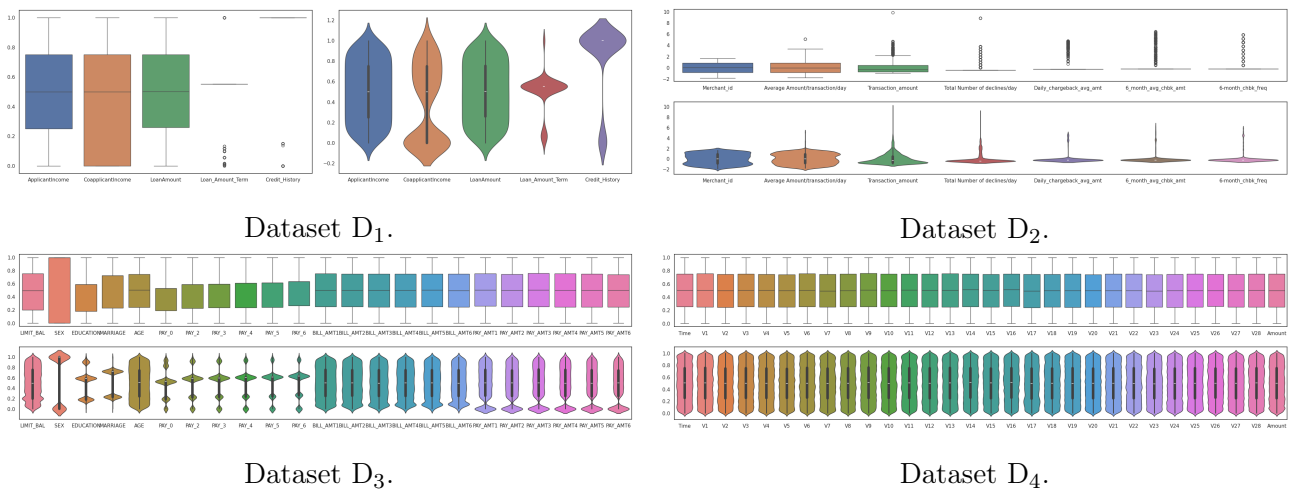


Figure 5: Box plots and violin plots of the normalized data in the datasets.

2.6. Imbalanced Data

In order to address the class imbalance present in the training data, the Synthetic Minority Over-sampling Technique (SMOTE) was employed.

This method surpasses traditional oversampling techniques, which merely duplicate existing minority class instances, by generating synthetic instances in the neighborhood of existing minority class samples.

SMOTE accomplishes this by identifying the k -nearest neighbors of every existing minority class samples x .

Then, it generates a synthetic instance at some intermediate point along the line segment joining x to one of its randomly chosen k -nearest neighbor, until the desired number of synthetic instances is reached. This approach effectively augments the minority class without altering the number of majority class instances, resulting in a more balanced dataset while preserving data diversity.

3. HybridIG-CSO

The key elements of the HybridIG-CSO model are feature selection and classification using Random Weight Network (RWN).

3.1. Feature selection

3.1.1. Information Gain (IG)

In the filter approach, features are evaluated independently based on statistical measures such as information gain. The information gain technique is a mathematical method, which measures the reduction in uncertainty, or entropy, about a target variable when a specific feature is known.

Mathematically, the information gain (IG) for a feature X with respect to a target variable Y can be calculated as:

$$IG(X|Y) = H(X) - H(X|Y) \quad (1)$$

where $H(Y)$ represents the entropy of the target variable Y before considering feature X , and $H(Y|X)$ represents the conditional entropy of Y given the values of feature X .

```
1 def InformationGain(attribute, target):
2     original_entropy = DistributionEntropy(target.value_counts())
3     attribute_entropy = AttributeEntropy(attribute, target)
4     return original_entropy - attribute_entropy
```

Additionally, the calculations for entropy $H(X)$ and conditional entropy $H(X|Y)$ are outlined as follows, where $P(x)$ is the proportion of instances with value x for feature X and $H(Y|X = x)$ is the entropy of Y for instances where feature X has value x :

$$H(X) = - \sum_{x \in X} P(x) \log_2(x) \quad (2)$$

```
1 def DistributionEntropy(fc):
2     d = fc/sum(fc)
3     return sum([-x * np.log2(x) for x in d[d!=0]])
```

$$H(X|Y) = - \sum_{x \in X} P(x) \sum_{y \in Y} P(x|y) \log_2(P(x|y)) \quad (3)$$

```
1 def AttributeEntropy(attribute, target):
2     # Attribute entropy for numerical attributes
3     if attribute.dtype != 'object':
4         values = np.sort(attribute.unique())
5         entropy_values = []
6         for i, value in enumerate(values):
7             left_mask = attribute <= value
8             right_mask = ~left_mask
9             entropy = left_mask.sum()/len(attribute) * DistributionEntropy(target[left_mask].
value_counts())
10            entropy += right_mask.sum()/len(attribute) * DistributionEntropy(target[right_mask].
value_counts())
11            entropy_values.append(entropy)
12        return np.min(entropy_values)
13    # Attribute entropy for categorical attributes
14    weights = attribute.value_counts()/sum(attribute.value_counts())
15    attribute_values = attribute.value_counts().index
16    entropy = 0
17    for i, value in enumerate(attribute_values):
18        mask = attribute == value
19        entropy += weights[i] * DistributionEntropy(target[mask].value_counts())
20    return entropy
```

The attribute entropy for continuous numerical attributes was determined by initially sorting numerical values, including class labels, followed by an exhaustive evaluation of potential cut-points to identify the optimal split maximizing information gain, equivalently minimizing attribute entropy.

3.1.2. Competitive Swarm Optimization (CSO)

The wrapper-based technique employed is the CSO, chosen for its ability to explore the complex search space of potential feature subsets. CSO is an algorithm rooted in the original PSO technique, devised to tackle the issue of premature convergence that often arises when applying PSO to complex search spaces containing numerous local optima^[6], by removing the influence of *gbest* and *pbest* associated with each particle. Consequently, CSO solely relies on pairwise competition between randomly selected particles to guide the search process. Hence, with the assumption of having k particles within the swarm (population), the CSO procedure initiates with

a population consisting of randomly initialized particles denoted as $P(t)$, where 't' signifies the generation. Every potential solution is depicted by one of the swarm's particles. Each particle can be considered as a point represented by a position X and a velocity V . During each iteration, the swarm $P(t)$ is divided into two equal and randomly selected groups. Subsequently, CSO selects two particles, one from each group, and initiates a competition between them. The winning particle is directly transferred to the next generation. The losing particle is updated based on the information derived from the winner and then included in the next generation. The positions and velocities of the particles that emerged as winners and losers in the i^{th} pairwise competition, relative to generation t , can be expressed as follows, where the winning particle's position is denoted as $X_{wi}(t)$ and its velocity $V_{wi}(t)$ and, similarly, the losing particle's position is $X_{li}(t)$ and its velocity $V_{li}(t)$:

$$V_{li}(t+1) = R_1(i,t)V_{li}(t) + R_2(i,t)(X_{wi}(t) - X_{li}(t)) + \varphi R_3(i,t)(\bar{X}_i(t) - X_{li}(t)) \quad (4)$$

$$X_{li}(t+1) = X_{li}(t) + V_{li}(t+1) \quad (5)$$

$R_1(i,t)$, $R_2(i,t)$ and $R_3(i,t)$ represent three vectors of randomly generated numbers sampled from the range $[0,1)$. $\bar{X}_i(t)$ denotes the average position of the pertinent particles. These pertinent particles could encompass the entire swarm of particles or a predefined group of neighboring particles. The parameter φ governs the extent to which $\bar{X}_i(t)$ influences the process.

3.2. Random Weight Network (RWN)

The fundamental architecture of the RWN architecture follows a fully connected architecture with a single hidden layer. Unlike conventional gradient descent methods that necessitate the adjustment of multiple parameters, such as learning rates and the number of training epochs, RWN simplifies the training process by focusing solely on the number of hidden neurons.

RWN begins by randomly initializing the input weights and hidden layer biases. Subsequently, it utilizes a dataset of N training samples to construct the hidden layer output matrix. The output weights are then determined through the application of the Moore-Penrose (MP) generalized inverse.

Given a dataset of N training samples, where each sample is represented as (x_i, t_i) , with $x_i \in \mathbb{R}^n$ and $t_i \in \mathbb{R}^m$, and considering the activation function $g(x)$, the weight vector $w_j \in \mathbb{R}^n$ connecting the j^{th} hidden neuron to the n input nodes, and the set of output weights values $\beta_j \in \mathbb{R}^m$ connecting the j^{th} hidden neurons to the m output nodes, compute the output weight matrix β [7], such that:

$$H\beta = T \quad (6)$$

where H is the hidden layer output matrix and T is the target matrix:

$$H = \begin{Bmatrix} g(w_1 \cdot x_1 + b_1) & \cdots & g(w_L \cdot x_1 + b_L) \\ \vdots & \cdots & \vdots \\ g(w_1 \cdot x_N + b_1) & \cdots & g(w_L \cdot x_N + b_L) \end{Bmatrix}_{N \times L} \quad \beta = \begin{bmatrix} \beta_1^T \\ \vdots \\ \beta_L^T \end{bmatrix}_{L \times m} \quad T = \begin{bmatrix} t_1^T \\ \vdots \\ t_N^T \end{bmatrix}_{N \times m} \quad (7)$$

The developed implementation utilized the sigmoid activation function in conjunction with a prediction threshold of 0.5. Predictions were classified as positive if they exceeded this threshold, and negative otherwise.

```

1 def sigmoid(x):
2     return 1 / (1 + np.exp(-x))
3
4 def compute_H(dataset, particle):
5     n_neurons = int(particle['n_neurons'], 2)
6     params = particle['params']
7     # Select only the weights and biases corresponding to the selected features
8     mask = np.array(list(particle['features'].values()))==1
9     weights = params[0][mask, :n_neurons]
10    bias = params[1][:n_neurons].flatten()
11    H = sigmoid(dataset @ weights + bias)
12    return np.array(H)
13
14 def RWN_model(particle, X, y):
15    X_train, y_train = prepare_data(particle, X, y)
16    # Calculate the hidden layer output matrix H
17    H = compute_H(X_train, particle)
18    # Output weights determined using the Moore-Penrose generalized inverse
19    return np.linalg.pinv(H) @ y_train

```

3.3. Methodology

The initial phase employs IG as a filter-based method to rank the features within the credit card dataset. Only the top-ranked features are retained and subsequently passed to the wrapper-based algorithm, CSO. The RWN algorithm serves as the learning model within this hybrid framework.

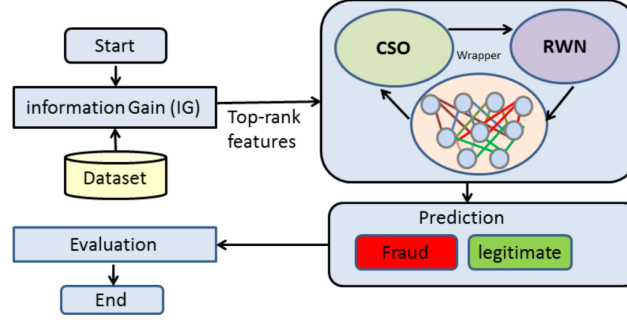


Figure 6: Flowchart of proposed method.

The CSO particle is encoded as a real vector encompassing the subsequent components:

- A set of binary flags indicating the inclusion or exclusion of corresponding features.
- A set of binary flags dictating the number of neurons in the hidden layer of the RWN.
- The RWN parameters, which encapsulate the values of input weights and hidden biases.

```

1 def create_particles(population, ig_attributes):
2     particles = []
3     for subset in population:
4         features = {}
5         for col in ig_attributes:
6             features[col] = 1 if col in subset else 0
7         n_neurons = ''.join(map(str, features.values()))
8         params = initialize_RWN_params(
9             np.array(list(features.values()), int(n_neurons, 2))
10        particles.append({
11            'features': features, 'n_neurons': n_neurons, 'params': params})
12    return particles

```

The fitness of each particle is evaluated to determine its effectiveness in contributing to the desired outcomes. This fitness calculation guides the optimization process. The fitness function is calculated as follows:

$$Fit = \alpha CLErr + \beta \frac{ft}{FT} + \gamma \frac{hd}{HD} \quad (8)$$

where $CLErr$ represents the error rate in classifying the RWN network, ft indicates the number of features selected, FT represents the overall count of features in the dataset, hd denotes the count of hidden neurons set by the optimizer and HD is the maximum allowable number of neurons in the RWN. The parameters α , β and γ manage the impact of weights, aiming to enhance the reduction rate of features, curtail RWN complexity, and diminish the quantity of chosen features.

```

1 def compute_fitness(particle, X_train, y_train, X, y, dataset_features, alpha, beta, gamma,
2     threshold):
3     # Compute the classification error rate
4     pred = compute_pred(particle, X_train, y_train, X, threshold)
5     error_rate = compute_error_rate(pred, y)
6     # Number of features identified
7     ft = np.sum(list(particle['features'].values()))
8     # Count of hidden neurons set by the optimizer
9     hd = int(particle['n_neurons'], 2)
10    # Maximum allowable number of neurons in the RWN
11    HD = 2*(len(particle['features'])+1) - 1
12    # Compute fitness
13    return alpha * error_rate + beta * (ft / dataset_features) + gamma * (hd / HD)

```


3.3.1. Attribute Ranking using IG Technique

In order to establish a suitable threshold, the standard deviation of IG values is calculated, a common practice for precise threshold determination^{[8][9]}. Attributes exceeding or meeting this threshold are retained, while those falling below it are discarded. This step ensures that only the most impactful attributes move forward.

```
1 ig_values = {}
2 for col in X_train.columns:
3     ig_values[col] = InformationGain(X_train[col], y_train)
4     print(f"{col}: {ig_values[col]}")
5 # Compute the standard deviation of the Information Gain (IG) values
6 threshold = np.std(list(ig_values.values()))
7 print(f"\nThreshold: {threshold}")
8 # Select attributes exceeding or meeting the threshold
9 ig_attributes = [col for col, ig in ig_values.items() if ig >= threshold]
10 print(f"Attributes exceeding or meeting the threshold: {ig_attributes}")
```

3.3.2. Wrapper CSO Algorithm with RWN Learning

The top-ranked attributes, identified through IG, are integrated into the wrapper CSO algorithm, with the RWN serving as the learning algorithm. The CSO's primary objective is to iteratively discover optimal subsets of features, achieved through a sequence of generations.

In the developed implementation, the RWN employs a pre-allocated maximum-sized zero weight matrix to enhance flexibility and computational efficiency. This strategic decision offers several advantages. Firstly, it accommodates the dynamic addition of new features. Secondly, by avoiding frequent memory reallocations, it significantly improves the algorithm's computational performance. Lastly, it provides a clear and intuitive representation of the feature-neuron relationship. During the update process of the CSO particle positions, the subset of features can change, necessitating a weight matrix with dimensions that can accommodate any possible feature subset. By initializing a matrix of dimensions equal to the maximum possible size and initially setting only weights corresponding to selected features using the uniform Xavier initialization, the RWN can dynamically adapt to any feature subset encountered by the CSO algorithm. Furthermore, the velocity update formula for the RWN weights within the CSO algorithm requires weight matrices of consistent dimensions for the velocity calculation. This is crucial as the update formula incorporates the current particle's weight matrix, the competitor's weight matrix and the mean weight matrix of the population. By defining a weight matrix of the maximum size and initializing only the relevant rows, we maintain compatibility with the CSO algorithm's velocity update requirements. This approach facilitates seamless updates and ensures that the optimization process remains consistent and effective, regardless of the changing feature subsets during the iterations.

```
1 def uniform_Xavier_initialization(input_size, output_size):
2     limit = np.sqrt(6 / (input_size + output_size))
3     return np.random.uniform(-limit, limit, (input_size, output_size))
4
5 def initialize_RWN_params(features, n_neurons):
6     params = list()
7     # Create a zero matrix of the max size
8     max_n_neurons = 2*(len(features)) - 1
9     W = np.zeros((len(features), max_n_neurons))
10    # Initialize only the selected portion with random weights
11    mask = features==1
12    W[mask, :n_neurons] = uniform_Xavier_initialization(np.sum(features), n_neurons)
13    # Create a zero vector of the max size
14    b = np.zeros((max_n_neurons, 1))
15    # Initialize only the selected portion with random weights
16    b[:n_neurons] = np.random.randn(n_neurons, 1)
17    params.append(W)
18    params.append(b)
19    return params
```

3.3.3. CSO Particle Initialization and Competition

The CSO Randomly initializes the number of individuals for the population, where candidate feature subsets are encoded as particles. Consequently, population is defined by enumerating all possible subsets of top-ranked features (combinations without repetition), setting a maximum size of 50, as suggested in the paper.

```
1 def population_generator(ig_attributes, max_size=50):
2     population = list(
3         chain.from_iterable(
```

```

4         combinations(ig_attributes, r) for r in range(1, len(ig_attributes)+1)
5     ))
6     if len(population) > max_size:
7         # Randomly select particles from the population within the max size limit
8         idxs = np.random.choice(len(population), max_size, False)
9         return [population[i] for i in idxs]
10    return population

```

The population is divided into two equal parts, each comprising $k/2$ particles. The ensuing pairwise competition determines winners and losers among particles. Winners advance to the next generation, while losers undergo updates before moved to the next generation. The subsequent stage entails training diverse RWNs using each particle, followed by the computation of the fitness value for each feature subset. Aims to refine the particle population iteratively.

The CSO algorithm's concepts of position and velocity were implemented by defining three positions and velocities for each particle: one for feature subset selection and two for the RWN configuration (one for weights and one for biases). These were initialized as follows:

```

1 velocities = [
2     [
3         np.random.randn(len(ig_attributes)),
4         np.random.randn(*particle['params'][0].shape),
5         np.random.randn(*particle['params'][1].shape)
6     ] for particle in particles
7 ]

```

The evaluation and update of particles were executed in parallel using Joblib, leveraging multiple processors to enhance computational efficiency and reduce execution time.

```

1 def competition(ig_attributes, particles, velocities, X_train, y_train, X_valid, y_valid,
2     dataset_features, alpha, beta, gamma, phi, threshold):
3     # Randomly shuffle the particles and velocities
4     particles, velocities = shuffle(particles, velocities)
5     majority_features, mean_params = compute_mean_particle(particles)
6     next_gen_particles = []
7     next_gen_vel = []
8     # Execute the evaluation and update of particles in parallel
9     results = Parallel(n_jobs=-1)(
10         delayed(evaluate_particles)(
11             ig_attributes,
12             particles[i], particles[i+1], velocities[i], velocities[i+1],
13             X_train, y_train, X_valid, y_valid,
14             dataset_features, alpha, beta, gamma, phi,
15             majority_features, mean_params,
16             threshold
17         ) for i in range(0, len(particles)-1, 2))
18     # Extract the results
19     for result in results:
20         next_gen_particles.extend(result[0])
21         next_gen_vel.extend(result[1])
22     # If there is an odd number of particles, add the last one to the winners list
23     if len(particles) % 2 != 0:
24         next_gen_particles.append(particles[-1])
25         next_gen_vel.append(velocities[-1])
26     return next_gen_particles, next_gen_vel

```

The update of the losing particle is carried out in accordance with the previously outlined CSO algorithm formulas:

```

1 def update_loser(ig_attributes, velocity, winner, loser, majority_features, mean_params, phi,
2     threshold):
3     # Matrix of vectors with randomly generated numbers sampled from the range [0,1)
4     R = np.random.rand(3)
5     # Compute the velocity update for feature selection in the loser particle
6     X_wi = np.array(list(winner['features'].values()))
7     X_li = np.array(list(loser['features'].values()))
8     vel_feature = R[0]*velocity[0] + R[1]*(X_wi - X_li) + R[2]*phi*(majority_features-X_li)
9     X_li_updated = (sigmoid(X_li + vel_feature) > threshold).astype(int)
10    features = {}
11    for i, col in enumerate(ig_attributes):
12        features[col] = X_li_updated[i]
13    n_neurons = ','.join(map(str, features.values()))
14    # Compute the velocity update for the weights in the RWN configuration
15    X_wi = winner['params'][0]
16    X_li = loser['params'][0]
17    vel_weights = R[0]*velocity[1] + R[1]*(X_wi - X_li) + R[2]*phi*(mean_params[0]-X_li)
18    weights = X_li + vel_weights

```

```

18 # Compute the velocity update for the biases in the RWN configuration
19 X_wi = winner['params'][1]
20 X_li = loser['params'][1]
21 vel_biases = R[0]*velocity[2] + R[1]*(X_wi - X_li) + R[2]*phi*(mean_params[1]-X_li)
22 biases = X_li + vel_biases
23 particle = {'features':features, 'n_neurons':n_neurons, 'params':[weights, biases]}
24 return particle, [vel_feature, vel_weights, vel_biases]

```

where the average position of the population's particles $\bar{X}_i(t)$ is determined by computing the majority vote for the selected subset of features and by calculating the mean value for each parameter within the weights and biases matrices.

```

1 def compute_mean_particle(particles):
2     # Compute the majority vote for selected features
3     data = np.array([list(particle['features'].values()) for particle in particles])
4     # Compare counts to half the number of rows to determine majority
5     majority_features = (np.sum(data, axis=0) > data.shape[0] / 2).astype(int)
6     # Compute the mean of the parameters for the particles
7     mean_weights = np.mean([particle['params'][0] for particle in particles], axis=0)
8     mean_biases = np.mean([particle['params'][1] for particle in particles], axis=0)
9     mean_params = [mean_weights, mean_biases]
10    return majority_features, mean_params

```

3.3.4. Iterative Refinement

The methodology persists iteratively until a predefined maximum iteration count is reached. Throughout this process, the wrapper CSO algorithm continues its pursuit of the best feature subset and corresponding RWN configuration.

3.4. Hyperparameters tuning

The hyperparameters α , β , γ and φ were optimized through a 10-fold cross-validation approach, by determining the combination that yields the highest average accuracy score over the 10 folds.

```

1 def hyperparameters_tuning(ig_attributes, X, y, dataset_features, categorical_cols, max_iter,
    threshold):
2     # SMOTE
3     X_train, y_train = perform_SMOTE(X, y, categorical_cols)
4     values = [1e-3, 0.1, 1, 10, 100]
5     # Provides train/test indices to split data in train/test sets
6     kf = KFold(n_splits=10, shuffle=True, random_state=seed)
7     # Define alpha search space
8     param_grid = [(alpha, 0.5, 0.5, 0.5) for alpha in values]
9     print('Tuning alpha ...')
10    best_score, best_params = perform_tuning(
11        kf, param_grid, -np.inf, None,
12        ig_attributes,
13        X, y, dataset_features, categorical_cols, max_iter, threshold)
14    # Define beta search space
15    param_grid = [(best_params[0], beta, 0.5, 0.5) for beta in values]
16    print('Tuning beta ...')
17    best_score, best_params = perform_tuning(
18        kf, param_grid, best_score, best_params,
19        ig_attributes,
20        X, y, dataset_features, categorical_cols, max_iter, threshold)
21    # Define gamma search space
22    param_grid = [(best_params[0], best_params[1], gamma, 0.5) for gamma in values]
23    print('Tuning gamma ...')
24    best_score, best_params = perform_tuning(
25        kf, param_grid, best_score, best_params,
26        ig_attributes,
27        X, y, dataset_features, categorical_cols, max_iter, threshold)
28    # Define phi search space
29    param_grid = [(best_params[0], best_params[1], best_params[2], phi) for phi in values]
30    print('Tuning phi ...')
31    best_score, best_params = perform_tuning(
32        kf, param_grid, best_score, best_params,
33        ig_attributes,
34        X, y, dataset_features, categorical_cols, max_iter, threshold)
35    print("Best parameters: ", best_params)
36    print("Best fitness score: ", best_score)
37    return best_params

```

4. Performance Evaluation

4.1. Performance Metrics

The evaluation of the involved models' performance is based on a series of metrics derived from the confusion matrix. These metrics utilize parameters such as true positives (TP), true negatives (TN), false positives (FP) and false negatives (FN):

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (9)$$

$$Sensitivity = \frac{TP}{TP + FN} \quad (10)$$

$$Specificity = \frac{TN}{TN + FP} \quad (11)$$

$$G - mean = \sqrt{Sensitivity \times Specificity} \quad (12)$$

$$Precision = \frac{TP}{TP + FP} \quad (13)$$

$$Recall = \frac{TP}{TP + FN} \quad (14)$$

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (15)$$

Finally, the AUC (Area Under the Curve), assessing the models' differentiation capability via the ROC curve.

4.2. Experiment I: Comparisons performance between HybridIG-CSO, RWN with filter approach, and RWN with CSO

In this experiment, the HybridIG-CSO method was evaluated by comparing it with three distinct techniques: the classical RWN, RWN with a filter-based approach (IG-RWN) and manually tuned CSO-RWN. The performance of the HybridIG-CSO method was assessed against the other approaches using six different metrics: Accuracy, Precision, Recall, AUC, F1 and G-mean.

Dataset	Algorithm	Accuracy	Precision	Recall	AUC	F1	G-Mean
D ₁	Classic-RWN	0.5393	0.5094	0.6429	0.5464	0.5684	0.5451
	IG-RWN	0.7528	0.9623	0.7183	0.7034	0.8226	0.6540
	CSO-RWN	0.6854	0.8302	0.6984	0.6512	0.7586	0.6261
	HybridIG-CSO	0.7528	0.9623	0.7183	0.7034	0.8226	0.6540
D ₂	Classic-RWN	0.9567	0.8765	0.8765	0.9251	0.8765	0.9239
	IG-RWN	0.9567	0.9383	0.8352	0.9495	0.8837	0.9494
	CSO-RWN	0.9827	0.9259	0.9740	0.9603	0.9494	0.9597
	HybridIG-CSO	0.9610	0.9506	0.8462	0.9569	0.8953	0.9569
D ₃	IG-RWN	0.7831	0.5494	0.5126	0.6998	0.5303	0.6834
	HybridIG-CSO	0.7793	0.5753	0.5044	0.7066	0.5375	0.6943
D ₄	IG-RWN	0.9784	0.8941	0.7525	0.9387	0.8172	0.9376
	HybridIG-CSO	0.9701	0.8824	0.6696	0.9288	0.7614	0.9276

Table 3: Performance experiment I on a separate test set.

Due to the limitations of the free usage plans on Colab and Kaggle, it was necessary to significantly reduce the number of features, as the initialization of the RWN exhausted all the available RAM. Consequently, it was not feasible to implement the classical RWN and the manually tuned CSO-RWN for the last two datasets, as these

require the use of the full dataset’s features. Regarding Dataset 4, downsampling was required. Specifically, the majority class was downsampled.

Dataset	Algorithm	Accuracy	Precision	Recall	AUC	F1	G-Mean
D ₁	Classic-RWN	0.4904	0.4984	0.4933	0.4912	0.4915	0.4857
	IG-RWN	0.7332	0.9540	0.6633	0.7350	0.7809	0.7004
	CSO-RWN	0.7790	0.8398	0.7513	0.7792	0.7913	0.7754
	HybridIG-CSO	0.7018	0.8960	0.6450	0.7023	0.7474	0.6713
D ₂	Classic-RWN	0.9720	0.9722	0.9721	0.9722	0.9721	0.9722
	IG-RWN	0.9730	0.9781	0.9683	0.9731	0.9732	0.9730
	CSO-RWN	0.9863	0.9892	0.9837	0.9865	0.9864	0.9865
	HybridIG-CSO	0.9707	0.9755	0.9665	0.9708	0.9709	0.9707
D ₃	IG-RWN	0.6980	0.5582	0.7749	0.6980	0.6490	0.6839
	HybridIG-CSO	0.6897	0.5490	0.7640	0.6897	0.6389	0.6752
D ₄	IG-RWN	0.9651	0.9502	0.9794	0.9651	0.9645	0.9649
	HybridIG-CSO	0.9212	0.9253	0.9339	0.9211	0.9268	0.9170

Table 4: Performance experiment I using 10-fold cross-validation.

4.3. Experiment II: Comparison with other classifiers

In this experimental evaluation, the effectiveness of the HybridIG-CSO is assessed in the context of fraud classification, by comparing its performance against other widely employed algorithms typically used as induction techniques in feature selection wrapper-based methods, namely Naïve Bayes (NB), Random Forest (RF) and Support Vector Machine (SVM).

Dataset	Classifier	Accuracy	Precision	Recall	AUC	F1	G-Mean
D ₁	NB	0.7528	0.9623	0.7183	0.7034	0.8226	0.6540
	RF	0.7528	0.9057	0.7385	0.7167	0.8136	0.6914
	SVM	0.7528	0.9623	0.7183	0.7034	0.8226	0.6540
	HybridIG-CSO	0.7528	0.9623	0.7183	0.7034	0.8226	0.6540
D ₂	NB	0.9113	0.7284	0.7564	0.8393	0.7421	0.8319
	RF	0.9762	0.9753	0.8977	0.9758	0.9349	0.9758
	SVM	0.9524	0.9877	0.7921	0.9663	0.8791	0.9660
	HybridIG-CSO	0.9610	0.9506	0.8462	0.9569	0.8953	0.9569
D ₃	NB	0.6769	0.6630	0.3734	0.6719	0.4777	0.6719
	RF	0.8047	0.4576	0.5781	0.6809	0.5109	0.6433
	SVM	0.7556	0.6022	0.4628	0.7009	0.5234	0.6939
	HybridIG-CSO	0.7793	0.5753	0.5044	0.7066	0.5375	0.6943
D ₄	NB	0.9886	0.8353	0.9467	0.9163	0.8875	0.9127
	RF	0.9886	0.8706	0.9136	0.9329	0.8916	0.9309
	SVM	0.9746	0.8941	0.7103	0.9366	0.7917	0.9357
	HybridIG-CSO	0.9701	0.8824	0.6696	0.9288	0.7614	0.9276

Table 5: Performance experiment II on a separate test set.

Dataset	Classifier	Accuracy	Precision	Recall	AUC	F1	G-Mean
D ₁	NB	0.7248	0.9492	0.6566	0.7267	0.7746	0.6907
	RF	0.8044	0.8334	0.7908	0.8050	0.8106	0.8040
	SVM	0.7248	0.9590	0.6550	0.7269	0.7765	0.6871
	HybridIG-CSO	0.7018	0.8960	0.6450	0.7023	0.7474	0.6713
D ₂	NB	0.8413	0.7125	0.9588	0.8411	0.8173	0.8311
	RF	0.9897	0.9955	0.9841	0.9898	0.9898	0.9898
	SVM	0.9673	0.9724	0.9628	0.9673	0.9675	0.9672
	HybridIG-CSO	0.9707	0.9755	0.9665	0.9708	0.9709	0.9707
D ₃	NB	0.6724	0.6657	0.6748	0.6724	0.6702	0.6723
	RF	0.8768	0.8525	0.8960	0.8768	0.8737	0.8764
	SVM	0.7017	0.5990	0.7538	0.7017	0.6675	0.6941
	HybridIG-CSO	0.6876	0.5455	0.7621	0.6876	0.6358	0.6727
D ₄	NB	0.9160	0.8358	0.9955	0.9160	0.9086	0.9124
	RF	0.9938	0.9906	0.9969	0.9938	0.9937	0.9937
	SVM	0.9619	0.9458	0.9773	0.9619	0.9613	0.9618
	HybridIG-CSO	0.9212	0.9253	0.9339	0.9211	0.9268	0.9170

Table 6: Performance experiment II using 10-fold cross-validation.

5. Conclusions

Before presenting the obtained results, it is important to acknowledge that accuracy, while commonly used metrics in classification tasks, is not a reliable index of the quality of the trained model in case of imbalanced datasets or when the importance of wrongly predicting positive-class samples is different from wrongly predicting negative-class samples. These metrics can give a misleading picture of the model’s performance, often favoring the majority class.

More suitable metrics in such scenarios are the AUC, F1 score and G-Mean, which provide a more balanced assessment by considering both precision and recall or by emphasizing the balance between sensitivity and specificity, thereby offering a more comprehensive evaluation of the model’s performance on imbalanced data. The performance achieved in these more robust metrics closely aligns with the results reported in the referenced paper, except for the Classic-RWN model in the first dataset in Experiment I, which exhibited notably superior performance.

The performance metrics obtained on a separate test set indicate that the HybridIG-CSO model outperforms the other models across all metrics and datasets, except for the second and fourth datasets. In these cases, the CSO-RWN and IG-RWN model, respectively, perform better, although the HybridIG-CSO demonstrates very close performance and superior Precision.

However, due to the limitations of the free usage plans on Colab and Kaggle, it was not feasible to implement the classical RWN and the manually tuned CSO-RWN for the last two datasets.

In the evaluation of performance using 10-fold cross-validation, which is more characterized by independence from the choice of samples, the CSO-RWN outperforms the HybridIG-CSO across all metrics. In contrast, for the last two datasets, IG-RWN emerges as the best-performing model among the two tested.

In the second experiment, the HybridIG-CSO outperforms the other classifiers in almost all metrics on the first and third datasets on a separate test set. In all other cases, including using 10-fold cross-validation, the ensemble learning method, Random Forest, which combines the output of multiple decision trees to reach a single result, emerges as the best classifier.

The hyperparameters α , β , γ and φ were optimized through a 10-fold cross-validation approach, by determining the combination that yields the highest average score over the 10 folds.

Obviously, the performance of HybridIG-CSO could be further optimized by expanding the search space for hyperparameters or optimizing with respect to other metrics such as F1, G-mean or AUC, increasing the number of iterations, testing different thresholds for filter-based feature selection using Information Gain, or testing a higher prediction threshold for the sigmoid function.

Ensemble learning with HybridIG-CSO models, using different configurations, could further enhance classification performance.

However, the most noteworthy aspect of this project is the efficiency of the method proposed in the paper.

With a significantly simplified training process compared to conventional gradient descent methods, it achieves performance, on a complex machine learning task as credit card fraud detection, that is almost comparable to an ensemble learning method like Random Forest, and in some cases, on a separate test set, even surpasses it. Particularly remarkable is the use of the hybrid approach for feature selection, combining Information Gain and Competitive Swarm Optimization, which enabled the model to be tested on large datasets, overcoming resource limitations, and outperforming the other models on a smaller dataset like the first one.

References

- [1] Enas Faisal, Hadeel Ahmad, and Rawan Zaghloul. Efficient credit card fraud detection using evolutionary hybrid feature selection and random weight networks. *International Journal of Data and Network Science*, 09 2024. doi: 10.5267/j.ijdns.2023.9.009.
- [2] Rejwan Bin Sulaiman, Vitaly Schetinin, and Paul Sant. Review of machine learning approach on credit card fraud detection. *Human-Centric Intelligent Systems*, 2, 05 2022. doi: 10.1007/s44230-022-00004-0.
- [3] Mohammad Masoud, Yousef Jaradat, Esraa Rababa, and Ahmad Manasrah. Turnover prediction using machine learning: Empirical study. *International Journal of Advances in Soft Computing & Its Applications*, 13(1), 2021.
- [4] Rafael Lima and Adriano Pereira. Feature selection approaches to fraud detection in e-payment systems. volume 278, pages 111–126, 02 2017. ISBN 978-3-319-53675-0. doi: 10.1007/978-3-319-53676-7_9.
- [5] Naoufal Rtayli and Nourddine Enneya. Enhanced credit card fraud detection based on svm-recursive feature elimination and hyper-parameters optimization. *Journal of Information Security and Applications*, 55: 102596, 12 2020. doi: 10.1016/j.jisa.2020.102596.
- [6] Ran Cheng and Yaochu Jin. A competitive swarm optimizer for large scale optimization. *IEEE transactions on cybernetics*, 45, 05 2014. doi: 10.1109/TCYB.2014.2322602.
- [7] Guang-Bin Huang, Qin-Yu Zhu, and Chee Siew. Extreme learning machine: A new learning scheme of feedforward neural networks. volume 2, pages 985 – 990 vol.2, 08 2004. ISBN 0-7803-8359-1. doi: 10.1109/IJCNN.2004.1380068.
- [8] J Roseline, GBSR Naidu, Samuthira Pandi Velusamy, S Rajasree, and Dr.N. Mageswari. Autonomous credit card fraud detection using machine learning approach. *Computers and Electrical Engineering*, 102:108132, 09 2022. doi: 10.1016/j.compeleceng.2022.108132.
- [9] Maria Prasetyowati, Nur Maulidevi, and Kridanto Surendro. Determining threshold value on information gain feature selection to increase speed and prediction accuracy of random forest. *Journal of Big Data*, 8, 06 2021. doi: 10.1186/s40537-021-00472-4.