

# A Quantum Beta Distributed Multi-Objective Particle Swarm Optimization Algorithm for Twitter Fake Accounts Detection

Ahlem Aboud <sup>a, b, c, \*</sup>, Nizar Rokbani <sup>b, c</sup>, Seyedali Mirjalili <sup>d, e</sup>, Amir Hussain <sup>f</sup>, Habib Chabchoub <sup>g</sup> and  
Adel M. Alimi <sup>b, h</sup>

<sup>a</sup> University of Sousse, ISITCom, 4011, Sousse, Tunisia.

<sup>b</sup> REGIM Lab: REsearch Groups in Intelligent Machines, University of Sfax, National Engineering School of Sfax (ENIS), BP 1173, Sfax, 3038, Tunisia.

<sup>c</sup> High Institute of Applied Science and technology of Sousse, University of Sousse, Tunisia.

<sup>d</sup> Centre for Artificial Intelligence Research and Optimisation, Torrens University Australia, Brisbane, Australia.

<sup>e</sup> Yonsei Frontier Lab, Yonsei University, Seoul, South Korea.

<sup>f</sup> Edinburgh Napier University, School of Computing, Edinburgh EH10 5DT, Scotland, U.K.

<sup>g</sup> College of Business, Al Ain University of Science and Technology, Abu Dhabi, United Arab Emirates.

<sup>h</sup> Department of Electrical and Electronic Engineering Science, Faculty of Engineering and the Built Environment, University of Johannesburg, South Africa.

-----  
\* Corresponding Author

E-mail Address: [ahlem.aboud@regim.usf.tn](mailto:ahlem.aboud@regim.usf.tn) (A. Aboud)

Contributing Authors: E-mail addresses: [nizar.rokbani@ieee.org](mailto:nizar.rokbani@ieee.org) (N. Rokbani), [ali.mirjalili@gmail.com](mailto:ali.mirjalili@gmail.com) (S. Mirjalili), [a.hussain@napier.ac.uk](mailto:a.hussain@napier.ac.uk) (A. Hussain), [habib.chabchoub@gmail.com](mailto:habib.chabchoub@gmail.com) (H. Chabchoub), [adel.alimi@ieee.org](mailto:adel.alimi@ieee.org) (A. M. Alimi).

## Abstract

Fake account detection is a topical issue when many Online Social Networks encounter several issues caused by the growing number of unethical online activities. This study presents a new Quantum Beta-behaved Multi-Objective Particle Swarm Optimization (QB-MOPSO) algorithm for machine learning based Twitter fake accounts detection. The proposed approach aims to improve the learning process of deep neural networks, random forest, through minimizing simultaneously the feature dimensionality and the classification error rate. The main contribution consists in proposing a quantum beta MOPSO to handle the training phase of neural and deep architectures. The QB-MOPSO is used to perform a multi-objective training of the random forest algorithm. The QB-MOPSO has two optimization profiles: the first one uses a quantum-behaved equation for improving the exploratory behaviour of PSO, while the second one uses a beta function to enhance PSO's exploitation. An extensive experimental study is carried out using two open Twitter datasets with 1982 and 928 accounts. The new proposal is a random forest QB-MOPSO. Results showed that random forest QB-MOPSO accuracy is about 99.19% and 97.52% accounts on datasets 1 and 2. Comparative analysis of the proposed architecture toward the original architecture showed that the use of QB-MOPSO for learning enhances the random forest algorithm which perform then the original ones.

**Keywords:** Feature Selection, Fake Account Detection, Quantum Beta Multi-Objective Particle Swarm Optimization, Machine Learning, Distributed System, Quantum Computing.

## 1. Introduction

Online Social Networks (OSNs) have become a crucial part of daily life. Social media and mobile devices are driving the growth of the World Wide Web. According to the digital report<sup>1</sup> published in January 2020, out of 7.75 billion people worldwide, there are 5.19 billion phone users, 4.54 billion Internet users, 3.8 billion active social media users, and 3.75 billion mobile social media users. The world's internet users spend an average of six hours online each day. In the last decade, many users have become addicted to the use of well-known online social networks like LinkedIn, Facebook, Twitter, YouTube, ...etc. It is not only for good habits like communication and sharing information, but a substantial part of OSN users is not human was presenting fake or bot accounts controlled by computers to gain popularity and promote business activities for financial gain. Several techniques are developed to manage the user profiling problem for capturing information about users and their interests which are called the User Data Discovery (UDD) model referred to the Knowledge Data Discovery (KDD) model. In this context, several approaches have developed and regrouped into three categories: explicit [1], implicit and hybrid user profiling techniques while the main issues covered the process of information retrieval and collection of the user's information [2]. Implicit user profiling approaches are referred to as static or factual profiling that provides the static process to analyse and collect static and predictable characteristics about users by filing some online forms. However, explicit approaches are referred to behavioral, adaptive, and ontological profiling that leads to the dynamic process of collecting future behaviors and learning about users based on several filtering techniques [1] such as Rule based filtering,

<sup>1</sup> <https://wearesocial.com/fr/blog/2020/01/digital-report-2020>

Collaborative filtering, and content-based filtering. Furthermore, hybrid approaches have combined the advantage of both explicit and implicit methods taking into consideration static and dynamic characteristics of the user profile to maintain the accuracy of temporal information.

During the last five years, a variety of approaches have been developed to manage user profiling problems not only with regards to data discovery but also for unhealthy activities detection like spam/ non-spam accounts [3], fake or bot accounts [4], fake followers [5], fake news [6], and fake engagement [7] using different Machine Learning Algorithms (MLAs) for classification purposes. The classification task involves five main steps namely, data collection, feature extraction, feature selection, classification, and prediction tasks. Therefore, the feature selection step has been considered a challenging problem for classification, clustering, time series prediction and regression tasks. This study focuses mainly on the Feature Selection (FS) problem [8] in the classification task to select a small subset of pertinent features to enhance the performance of machine learning models in terms of the best accuracy, and interpretability to minimise the computational time [9]. Thus, the Twitter online social network has attracted many researchers due to the severity of Twitter's social spambots problem [10] and the availability of public datasets which were easier to find than those of other OSNs.

The main motivation of this contribution is to propose a novel Quantum Beta-behaved Multi-Objective Particle Swarm Optimization algorithm (QB-MOPSO) to enhance the training process of machine learning algorithms based on a best subset of selected features for Twitter fake accounts detection. The proposed QB-MOPSO algorithm aims to examine the process of selecting relevant features that minimize the classification error rate. The optimization process of the QB-MOPSO algorithm has developed based on the Revised Quantum-behaved Particle Swarm Optimization (RQPSO) [11], and the Gaussian Quantum-behaved Particle Swarm Optimization algorithm (GAQPSO) [12] as well as the use of beta function provided by Alimi [13] in 2003. The QB-MOPSO algorithm presents two optimization profiles. In the first one, all particles are subject to quantum PSO approaches (RQPSO and GAQPSO) with random uniform and Gaussian distributions to better explore the search space. The second one is for exploitation enhancement using a beta function with three data distribution shapes namely; Gaussian, linear decrease, and exponential forms. The quantum and beta-behaved rules aim to assume a higher level of convergence toward the global best solutions. The proposed QB-MOPSO algorithm starts with random initialisation of  $N$  particles. Each particle  $P$  is a potential solution performed in the search space. The dynamic switching phases are assumed by the two optimization profiles which are symmetric about the mean personal best position ( $mbest$ ). A particle  $P$  is optimized for exploration phase, if the current best position ( $pbest$ ) is less than the mean best position ( $mbest$ ). Otherwise, it is considered for the exploitation phase. At each iteration, the update rules of particle position are as follows:

- **Exploration phase:** the particle position has been updated using the quantum equation in (RQPSO and GAQPSO).
- **Exploitation phase:** the particle position has been updated using the beta function.

The proposed application of the proposed QB-MOPSO approach for identifying fake accounts is denoted by the Neuro-QB-MOPSO method. To select pertinent features, a primordial step is added to the QB-MOPSO algorithm and named position binary encoding based on the sigmoid function. Only bits with the value of “1” are considered as selected features and used to train and test the classification model. For decision making process, the best subset of multiple

pertinent features is determined based on the compromise solution that minimize the classification error rate and determined using the nearest non-dominated solution to the utopian point.

The rest of this paper is resumed as follows: Section 2 presents the definition of feature selection problem and an overview of the existing approaches-based features selection techniques for fake accounts detection. Section 3, summarize the existing Quantum-behaved PSO Methods. The proposed Quantum Beta-behaved Multi-Objective Particle Swarm Optimization (QB-MOPSO) algorithm has detailed in Section 4. Section 5, presents the complexity analysis of the proposed QB-MOPSO approach. Section 6, details the Neuro-QB-MOPSO architecture for fake account detection. The preliminary of the experimental study and the comparative results are discussed in Section 7. Finally, Section 8 concludes the paper and suggests the future work.

## 2. State-of-the-Art

This section presents the definition of the Feature Selection (FS) problem and the existing approaches-based FS techniques for fake account detection.

### 2.1 Feature Selection Problem Statement

The Feature Selection (FS) has defined as a minimization Multi-Objective Optimization Problem (MOP) [9] to select a subset of relevant features to enhance the accuracy of MLAs for the classification task. The mathematical definition of FS problem [9] has presented in Equation (1). Let consider  $n$  data points  $X = \{x_i\}_1^n$  that presents the input dataset. Each sample  $x_i$  has d-dimensional features  $\{x_1, x_2, \dots, x_d\}$ . Two minimization objective functions are considered for the features dimensionality reduction  $F_1(x)$  in Equation (2) and the classification error rate  $F_2(x)$  in Equation (3).

$$\text{Minimize } F(x) = \begin{cases} F_1(X): \text{features dimensionality} \\ F_2(X): \text{classification error} \end{cases} \quad (1)$$

$$F_1(x) = \alpha * \frac{\#features}{\#All \ features} + (1 - \alpha) * \frac{Error\_Rate}{Error\_All} \quad (2)$$

where  $\alpha$  is a constant value  $\alpha \in [0,1]$ ,  $\#features$ : are the dimensionality of selected features,  $\#All \ Features$ : is the total number of original features.  $Error\_Rate$ : is the classification error rate of selected features.  $Error\_All$ : is the classification error rate using all features.

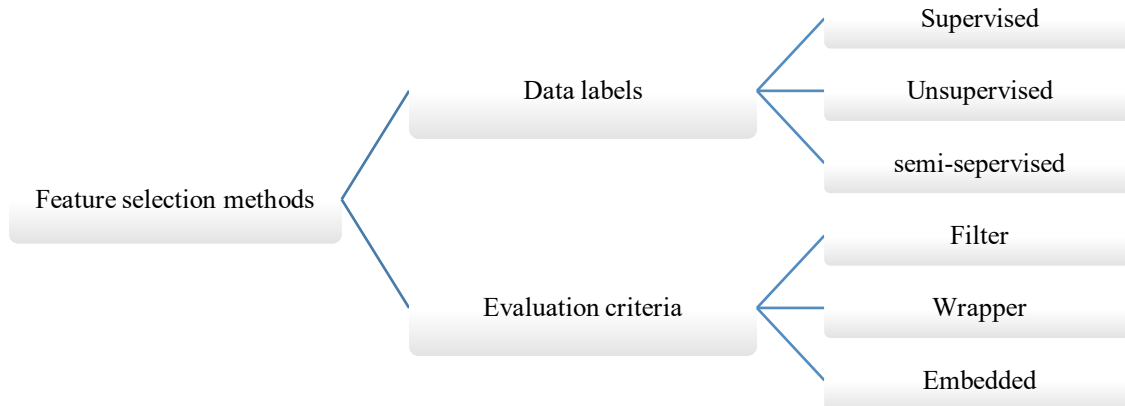
$$F_2(x) = \frac{FP+FN}{TP+TN+FP+FN} \quad (3)$$

where FP is the False Positive, FN is the False Negative, TP is the true positive, and TN is the true negative.

## 2.2 Existing Approaches-based Feature Selection Techniques for Fake Account Detection

On OSNs, several users aim to gain popularity not only by sharing healthy information about a specific domain of interest but also by introducing malicious activities, such as posting fake links and news. In 2016, the annual web traffic report<sup>2</sup> stated that more than 16.7 billion web visits to 100,000 randomly-selected web sites had analysed and detected more than 51.8% of bot users. In 2018, the industry report<sup>3</sup> announced 42.2% of all internet traffic was not human. Nevertheless, the increase of fake accounts generation has attributed to several extreme situations, including elections, Black Friday, COVID-19 pandemic and many other national or international events, activities, and diseases. A high numbers of user profiling techniques have developed and aimed to address different issues on OSNs like user interest detection, sentiment analysis, spam detection and fake account detection. Regarding to the availability of the online data of Twitter OSN, a high number of research studies have been publishing to addresses the problem of identifying Twitter fake accounts as a step toward fake news detection.

Online user profiles have included personal information's, shared content, links, and social interaction relationships. However, many online contents are shared to attract many users and a credible content is one of the most basic criteria for trusting web users. However, due to the widespread availability and usability of several features, many automatic programs (aka. bot accounts) are developed simulating human behavior. Such accounts are considered fake and impose malicious activities on social media platforms [4] and automatically controlled by a computer system to spread harmful activities [14]. The absence of a picture profile or online activities is not a very reliable method for detecting fake accounts. Therefore, it is very difficult to distinguish visually between bots and human accounts on OSNs. This has made fake account detection a challenging problem and has attracted several researchers.



**Figure 1.** Classification of feature selection methods

As resumed in Table 1, there are a wide range of techniques in the literature have been developing for feature selection to determine the most effective characteristic of a fake user. However, many features selection approaches

<sup>2</sup> <https://www.imperva.com/blog/bot-traffic-report-2016/>

<sup>3</sup> <https://www.globaldots.com/bad-bot-report-2018>

have proposed to enhance the classification task and presented in Figure 1. According to data labels three categories have proposed including; supervised, unsupervised and semi-supervised techniques. The main difference between the three categories has investigated respectively by the presence, the absence, and the existence of a small portion of labelled data [8]. The input labelled data makes supervised methods more specific for the decision making. Compared with unsupervised methods, the supervised approaches have produced a high accuracy, but the human intervention for the supervised learning need a high computational cost, and cannot be useful for the real-time data. Furthermore, the input unlabelled data makes a low complexity for unsupervised methods and the labels are determined automatically by the machine which are very useful for real-time data.

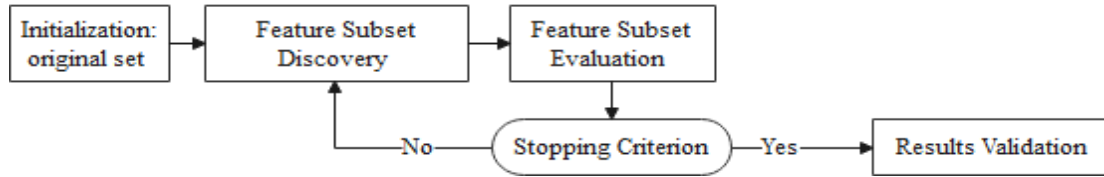
However, three categories of approaches are considered based on evaluation criteria including; filter, wrapper, and embedded methods. There have been both advantages and disadvantages to feature selection methods, depending on factors such as computational cost, speed, the dimension of the data, criteria for selecting features, and machine learning algorithms. Filter methods have been characterised by a high speed of treatment, a low computational cost, and well designed for a high dimensional data, however the use of statistical criteria does not guarantee the best subset of selected features [15]. The wrapper method has included a learning algorithm to determine the accuracy of the selected features [16], and to guarantee a better result compared with filter methods, but it has not performed with high dimensional data. Embedded techniques were the hybridization of both filter and wrapper methods. In filter methods, a statistical criterion has used for features dimensionality reduction and for wrapper methods the learning algorithm has considered to determine the best subset of features leading to the high classification accuracy [16].

In 2020, Rostami and Karbasi [17] used the Minimum Redundancy –Maximum Relevance algorithm (mRMR) [18] to identify the relevant subset of features with less redundancy. However, the previous feature selection techniques [4], [19] examined the best feature set based on the highest relation to the target class without taking into consideration the issue of independence and redundancy between selected features [20]. Ahmed and Abulaish [19], developed a generic statistical approach for spam detection-based Twitter and Facebook datasets. Azab *et al.* [4], have used the GAIN univariate algorithm for feature selection to determine the most effective subset of features that enhance the classification performance instead of using all features. In the most of cases, the use of statistical criterion does not guarantee the best subset of selected features [15].

Davis *et al.* [21], developed the BotOrNot platform using the Random Forest classifier as a black box approach for feature dimensionality reduction, and aims to evaluate whether a Twitter account is controlled by a human or machine. 1K features are extracted from the interaction patterns and the content. All collected features are regrouped into six classes of network features, user, friends, temporal, content, and sentiment features. Cresci *et al.* [22], proposed a Digital DNA model to predict online user behaviors such as new content, following or replying to other users. Yang *et al.* [23], have presented an empirical analysis of profile-based feature evasion tactics and content-based feature evasion tactics. Miller *et al.* [24], introduced a clustering model for anomaly detection. Moreover, different approaches have proposed to examine the stability of selected features by computing the similarity of the subset [25] or the use of machine learning algorithms to calculated the accuracy of the model using only the selected feature set [17].

Nevertheless, a variety of population-based approaches have designed for linear static and dynamic multi-objective optimization problems as well as for solving a set of complex problems involving at least three objective functions [26]–[32]. A set of evolutionary-based approaches like Genetic Algorithm (GA) [33], Particle Swarm Optimization

(PSO) [34], Genetic Programming (GP) [35], and Ant Colony Optimization (ACO) [36] are used for feature selection problem. Figure 2, details the iterative steps of feature selection process which are; initialization, feature subset discovery, feature subset evaluation and results validation. For feature selection methods, the key factors are the search techniques and evaluation criteria.



**Figure 2.** The iterative steps of the feature selection process

In 2015, Xue *et al.* [37] have published a survey paper to review the existing contributions based on population-based algorithms for solving single and multi-objective feature selection problem to minimise two objective functions (i) number of features and (ii) classification error rate. Search techniques, evaluation criteria, and the number of fitness functions are the three main concepts of evolutionary approaches. First, greedy search algorithms such as sequential forward selection (SFS) [38], and sequential backward selection (SBS) [39] are the well-known heuristic search techniques for feature selection. The main disadvantage of these methods is the ‘nesting effects’, so removed or selected features cannot be used for later testing. In addition, evolutionary techniques-based GA [40], GP [41], PSO [42], [43], and ACO [44] have been considered to determine which non-dominated solution provides the best trade-off between the number of features and the classification accuracy. Most of the existing feature selection methods suffer from the issues of high computational time/cost and stagnation in local optimum.

In contrast, few works use PSO for fake account detection, where the PSO algorithm is only used to optimize the parameters of both the logistic regression model [45], and the Q-learning method [46]. To sum up, the most existing contributions for fake account detection are done for minimising the number of the selected features taken into consideration the stability issue of the subset. The main goal was to enhance the computational time and maximise the accuracy rate of machine learning algorithms. The proposal consists in strongly exploring by a swarm having a quantum behavior before switching to a more stable behavior. This contribution has devoted to detect fake accounts on Twitter based on a new PSO-based approach whose role is to improve the self-learning of the deep neural detection system. Besides, several existing quantum-behaved PSO methods are reported in the next section.

**Table 1.** Existing methods for fake accounts detection on OSNs.

References	FS techniques	Tested Classifiers	Nb. Selected features	Accuracy	Datasets
Rostami and Karbasi [17], 2020	Minimum Redundancy Maximum Relevance algorithm [18].	10-fold cross validation using: (Random Forest, Naïve Bayes, SVM)	Test set 1: 8 Test set 2: 7 numerical	Best classifier: SVM Test set 1: 98% Test set 2: 97.1%	-Two Twitter datasets of Cresci <i>et al.</i> [10]
Cresci <i>et al.</i> [22], 2016	Digital DNA inspired by the biological DNA to model online user behaviors	ten-fold cross validation using: Bayes Net classifier	14 generic statistical features	Test set 1: 97.6% Test set 2: 92.9%	- <u>Dataset 1</u> : political - <u>Dataset 2</u> : Amazon
Davis <i>et al.</i> [21], 2016	Compute the bot-likelihood score using MLAs.	Ten-fold cross-validation using: Random Forest	1000 numerical feature values	95% AUC (Area Under ROC Curve).	-Dataset of 15k manually verified social bots and 16k legitimate accounts.
Azab <i>et al.</i> [4], 2016	GAIN univariate algorithm	5-fold cross validation using: (Random Forest, Decision Tree, Naïve Bayes, Neural Network, SVM)	7 numerical feature values	F-Measure (%) using: -RF:82.7, DT: 85.03, NB: 85.36, NN: 84.87, and SVM: 85.06	-Dataset of Twitter accounts collected by “the Fake project”
Miller <i>et al.</i> [24], 2014	Anomaly detection approach based on clustering model is built on normal twitter users with all outliers being treated as spam	Clustering algorithms (StreamKM++, DenStream, Combined)	126 numerical feature values	Accuracy using: - StreamKM++: 93.93% - DenStream: 97.11% - Combined: 98%	Dataset with 3239 user accounts including sample tweet (training set: 1587, test set: 1652)
Yang <i>et al.</i> [23], 2013	empirical analysis profile- based feature evasion tactics and content-based feature evasion tactics	10-fold cross validation using: (Random Forest, Decision Tree, Bayes Net, and Decorate)	25 numerical feature values	Best F1 Measure using Dataset I: RF :90%, Dataset II: RF :94.7%	- <u>Dataset I</u> : 20,000 accounts spam tweets, - <u>Dataset II</u> : 35,000 Twitter accounts
Ahmed and Abulaish [19], 2013	Generic statistical approach	Naïve Bayes, Jrip, and J48	14 generic statistical features	<u>Combined datasets</u> : detection rate (DR): 95.7%, false positive (FPR): 4.8% <u>Facebook dataset</u> : <u>DR</u> :96.4 %, FPR: 8.9%, <u>Twitter dataset</u> DR: 97.6%, FPR: 7.5%	Facebook and Twitter datasets



### 2.3 Existing Quantum-behaved PSO Methods

In 2004, Sun *et al.* [47] introduced quantum computing into the standard PSO algorithm. Quantum behaved PSO (QPSO) outperforms traditional PSO [34] with fewer control parameters and assumes a high level of convergence during the optimization process. So, instead of using a uniform stochastic distribution of particles' positions and velocity as in the original PSO algorithm. The quantum state of each particle is depicted by the wave function  $\Psi(x, t)$ ,  $\forall \lim_{x \rightarrow \pm\infty} \Psi(x) = 0$ . In quantum 3-dimensional time-space, the particle position in a point  $(x, y, z)$  is measured based on the probability density function  $|\Psi(x)|^2$  satisfying the normalization condition in Equation (4).

$$\int_{-\infty}^{+\infty} |\Psi(x)|^2 dx dy dz = \int_{-\infty}^{+\infty} Q dx dy dz = 1 \quad (4)$$

$$\text{Subject to: } \Psi(x) = \frac{1}{\sqrt{L}} \exp(-\|p - x\|/L)$$

where  $i\hbar \frac{\partial}{\partial t} \Psi(\vec{x}, t) = \hat{H} \Psi(\vec{x}, t)$  is the time-dependent Schrodinger equation,  $\hat{H}$  is the Hamiltonian operator,  $p$  is the center of potential and  $L$  is the vital parameter for creativity or imagination of the particle and computed at each iteration using the following equation  $L(t+1) = 2 \times \alpha \times |p - x(t)|$  and  $\lim_{t \rightarrow \infty} L(t) = 0$ . Sun *et al.* [47] considered the PSO algorithm as a quantum system where the quantum state is assumed using the wave function. The state of each particle at the time  $t$  is determined using the time-dependent Schrodinger equation  $|\Psi(x, t)|^2 = 1/L \exp(-2\|y\|/L)$ , where  $y = \pm \frac{L}{2} \ln(1/u)$ ,  $u = \text{rand}(0,1)$  and  $L$  is developed in time  $L = L(t) = 1/g |x_{id}(t) - p|$  and  $g > \ln\sqrt{2}$ .

#### a. Standard Quantum PSO (QPSO)

The first quantum-behaved PSO (QPSO) [47] is obtained through stochastic simulation of Monte Carlo measurement, when the particle position  $X(t)$  is given by:  $X(t) = p \pm \frac{L}{2} \ln(1/u)$ , with  $L_{(t+1)} = 2 \times \beta \times |p - X(t)|$  and the update equation  $x_i(t+1)$  of the particle  $i$  is presented in Equation (5).

$$x_i(t+1) = pbest_i^t \mp \beta \times |pbest_i^t - x_i^t| \times \ln\left(\frac{1}{u}\right) \quad (5)$$

The  $\beta$  parameter of QPSO is the contraction expansion factor of the algorithm (positive real number) is also called "Creativity" or "Imagination" of the particle. The update rules are dependant to the personal best position ( $pbest$ ) affected by a random uniform distribution of the parameter  $u$  uniformly distributed between 0 and 1.

To prevent the premature convergence, which a phenomenon that prevents an algorithm from finding an accurate estimation of the global optimum in meta-heuristics, a few QPSO improvements have been developed in the literature, including the Quantum Delta-Potential-Well-based Particle Swarm Optimization (QDPSO) algorithm [47], the Revised QPSO (RQPSO) [11], and the Gaussian Quantum-behaved PSO (GAQPSO) [12]. More details are presented in the next subsection.

### b. Quantum Delta-Potential-Well-based Particle Swarm Optimization (QDPSO)

In [47], Sun *et al.* have assumed that a quantum particle moves through a Delta potential well with a probability  $Z > 0.5$ . The particle is moved in a limited search space with respect to  $Z$ , otherwise it would appear out with a probability of  $1 - Z$ . The QDPSO algorithm benefits from the local attractor ( $La$ ) instead of the simple use of  $pbest$  position to guarantee convergence. In addition, the  $\beta$  coefficient is a positive real number set to:  $\beta = \frac{1}{g}, \forall g > \ln\sqrt{2}$ , which balances local and global searching ability. The update of the particle position in QDPSO algorithm is done using the Monte Carlo method in Equation (6).

$$x_i(t+1) = \begin{cases} \text{If } rand(0,1) \geq 0.5 \text{ then:} \\ La_i + L * (\ln(1/u)) \\ \text{Else:} \\ La_i - L * (\ln(1/u)) \end{cases} \quad (6)$$

where,

- $La_i = (\varphi_1 * pbest_i(t) + \varphi_2 * gbest(t)) / (\varphi_1 + \varphi_2)$
- $L = \beta * abs(La_i - x_i(t))$
- $u, \varphi_1$ , and  $\varphi_2 = rand(0,1)$
- $\beta = 1/g, \forall g > \ln\sqrt{2}$
- $T =$  shows the maximum number of iterations, and
- $gbest =$  is the global best position

### c. Revised Quantum PSO (RQPSO)

The global search ability of the QPSO system is denoted by the Revised Quantum PSO (RQPSO) [11]. The main difference between QDPSO and RQPSO is in the use of Mean Best Position ( $mbest$ ), which is denoted by the Mainstream Thought Point presenting the center-of-gravity global best position as presented in Equation (7).

$$mbest = \frac{1}{N} \sum_{i=1}^N pbest_i = \frac{1}{N} \sum_{i=1}^N pbest_{i1}, \dots, \frac{1}{N} \sum_{i=1}^N pbest_{in} \quad (7)$$

where,

- $N$  : indicates the size of the swarm, and.
- $mbest$  : represents mean global  $pbest$  position among all particles.

However, the equation to update particle position is modified by the parameter  $L$  which is equal to  $\beta * abs(mbest - x_i(t))$ ,  $\varphi_1$  and  $\varphi_2$  are two random parameters uniformly distributed between 0 and 1. The modified equation is presented in Equation (8).

$$x_i(t+1) = \begin{cases} \text{If } rand(0,1) \geq 0.5 \text{ then:} \\ La_i + L * (\ln(1/u)) \\ \text{Else:} \\ La_i - L * (\ln(1/u)) \end{cases} \quad (8)$$

where,

- $La_i = (\varphi_1 * pbest_i(t) + (1 - \varphi_2) * gbest(t)) / (\varphi_1 + \varphi_2)$ ,
- $L = \beta * abs(mbest - x_i(t))$ ,
- $u, \varphi_1$ , and  $\varphi_2 = rand(0,1)$ , and
- $\beta = 0.5 + 0.5 * (T - t) / T$ .

#### d. The Gaussian QPSO algorithm (GAQPSO)

The improved variant of QPSO system is denoted by the Gaussian QPSO algorithm (GAQPSO) [12], where the position is updated through a Gaussian distribution. The GAQPSO algorithm is developed to deal with the issue of trapping in the local optimum. In this case, the main modification between RQPSO and GAQPSO is in the random parameters  $u$ ,  $\varphi_1$ , and  $\varphi_2$  which are modified to follow a Gaussian probability distribution with zero mean and unit variance. The update of particle position in GAQPSO is done using Equation (9).

$$x_i(t+1) = \begin{cases} \text{If } rand(0,1) \geq 0.5 \text{ then:} \\ La_i + L * (\ln(1/u)) \\ \text{Else:} \\ La_i - L * (\ln(1/u)) \end{cases} \quad (9)$$

where,

- $La_i = (\varphi_1 * pbest_i(t) + (1 - \varphi_2) * gbest(t)) / (\varphi_1 + \varphi_2)$ ,
- $L = \beta * abs(mbest - x_i(t))$ ,
- $u, \varphi_1$ , and  $\varphi_2 = abs(N(0,1))$ ,
- $\beta = 0.5 + 0.5 * (T - t) / T$ ,
- $T$  = is the maximum number of iterations, and
- $gbest$  = indicates the global best position.

### 3. The Proposed Neuro- QB-MOPSO Architecture

This section presents two main parts:

- ✓ First a new Quantum Beta-behaved Multi-Objective Particle Swarm Optimization Algorithm (QB-MOPSO) is proposed
- ✓ Second, the new QB-MOPSO algorithm is applied to enhance the training process of neural architectures.

#### 3.1 The new QB-MOPSO Algorithm

The proposed QB-MOPSO algorithm benefits from two optimization profiles for exploration and exploitation phases based on a hybridisation of dynamic switching behaviors of two variants of quantum-behaved PSO approaches

(GAQPSO, and RQPSO) for the first exploration phase, and three types of parameters configuration of beta function using Gaussian, linear decrease, and exponential data distributions for the second exploitation phase.

The main different between the proposed system and the previous variants of quantum PSO has presented in position update rule and detailed in the following points:

- In the standard PSO algorithm [34] particles positions are updated according to the trajectory in Newtonian mechanics with a stochastic data distribution with respect to the current position and the new velocity computed using three components namely, the inertia weight, the personal best position, and the global best solutions in the population.
- The quantum variant of PSO (QDPSO) [47] aims to modify the update rule of the particle position using the quantum mechanics based on a time-dependent Schrodinger equation (see Equation 4) when the trajectory of particles has followed the wave function and the Delta potential well function with a probability equal to 0.5 (see Equation 6).
- In the revised quantum PSO algorithm (RQPSO) [11] the update rule of the particle position is like the QDPSO algorithm and the main difference has been investigated in the use of the center-of-gravity global best position denoted by the Mean Best Position (*mbest*).
- The Gaussian quantum PSO algorithm (GAQPSO) [12] has modified the update rules using the gaussian distribution form to update the personnel best solutions which is modified according to the normal distribution with zero mean and one standard deviation value.
- The proposed QB-MOPSO algorithm aims to hybridize both Newtonian mechanics and quantum mechanics to modify the rule to update the particle positions using two optimization profiles as presented in Equation (10). Inspired from the study of Sun *et al.* [47] which have proved that the quantum PSO algorithm has performed the traditional PSO algorithm with a high level of convergence leading to the best exploration of the search space. Moreover, the standard trajectory of particle position has been modified from a stochastic uniform distribution to a new data distribution with respect to the Beta function (Gaussian, linear decrease, and exponential). So, a particle  $P_i$  is performed in a quantum exploration profile, if their personnel best position (pbest ) is dominate or equal to the mean personal best position (*mbest*) ( $pbest \geq mbest$ ) otherwise it was considered for the second beta exploitation profile. The quantum and beta behaved rules aim to assume a higher level of convergence toward the global best solutions.

$$X_i(t+1) = \begin{cases} \rightarrow \textbf{Exploration Profile} \\ \text{IF } pbest_i(t) \geq mbest(t) \text{ then:} \\ \quad X_{t+1} = \text{update position using quantum} \\ \quad \quad \text{equation in RQPSO or GAQPSO} \\ \text{Else} \\ \rightarrow \textbf{Exploitation Profile} \\ \quad X_{t+1} = \text{update position using beta function} \end{cases} \quad (10)$$

The two dynamic switching behaviors are detailed as follows:

- **Behavior 1: exploration phase using quantum-behaved PSO**

For the exploration phase, all particles' positions are updated using the same equation in RQPSO and GAQPSO as explained in Equations (8) and (9) respectively.

- **Behavior 2: exploitation phase using beta-behaved PSO**

For the exploitation phase, the particle positions are updated using Equation (11) including the use of the beta function.

$$X_i(t+1) = X_t + V_{t+1} \quad (11)$$

where,  $V_{t+1}$  is the velocity of the particle has followed different distribution shapes according to the beta function, as presented in Equation (12).

$$V_{t+1} = V_t + \beta(x)(P_{best}(t) - X_t) + \beta(x)(g_{best}(t) - X_t) \quad (12)$$

where;  $g_{best}$  is the global best solution for all neighbours in swarm and  $p_{best}$  is the best local experience of each particle. Both  $g_{best}$  and  $p_{best}$  are used to affect the updated position of each particle at each iteration ( $t$ ).  $\beta(x)$  is the Beta function proposed by Alimi [13], and presented in Equation (13) and  $m_{best}$ , is computed using the Equation (7).

The beta function first introduced as a neural activation function [13] and demonstrated its ability to generate rich and flexible shapes (asymmetry, linearity, etc.). Also, the beta function has been adapted for different data distributions for feed-forward Neural Network (NN) [13] and investigated for Dynamic MOP and Many-Objective Optimization Problem [29]. According to the different configuration of both properties of  $p$  and  $q$ , three different shapes are considered in this study: the beta function with Gaussian, linear decrease, and exponential distributions as presented in the following Figure 3.

$$\beta(x) = \beta(x_0, x_1, p, q)(x) = \begin{cases} \text{if } x \in [x_0, x_1] \text{ then:} \\ \left( \frac{x-x_0}{x_c-x_0} \right)^p \left( \frac{x_1-x}{x_1-x_c} \right)^q \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

where,  $p, q, x_0, x_1$  are real numbers and  $x_c$  is the beta center defined in Equation (14).

$$x_c = \frac{px_1 + qx_0}{p+q} \quad (14)$$

where,

- $x_c$ : is the beta center point,
- $x_0$ , and  $x_1$  : are real numbers of the beta function in Equation (13), and
- $p$  and  $q$ : are the control properties of the beta function in Equation (13).

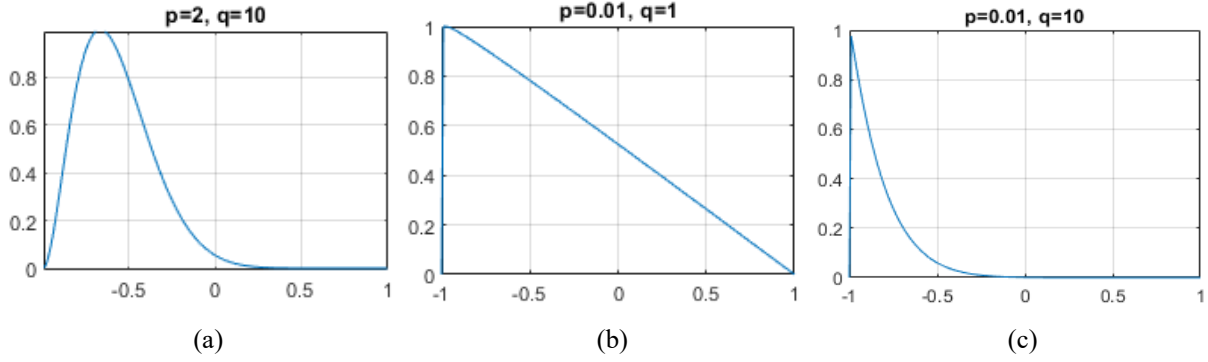
The multi-dimensional Beta function defined in Equation (15) presenting the product of  $m$  one-dimensional Beta function.

$$\beta(x) = \prod_{k=1}^m \beta(x_k, p_k, q_k, x_{0,k}, x_{1,k}) \quad (15)$$

where,

- $\prod_{k=1}^m \beta(x_k, p_k, q_k, x_{0,k}, x_{1,k})$  is the product of  $m$  one-dimensional Beta function in Equation (13),
- $m$ : is the dimension of the search space.

Based on different configurations of the parameters  $p$  and  $q$ , Figure 3 illustrates three examples of shapes that can be generated from the beta-function in Equation (13). According to Figure 3 (a), the beta-function has been assumed to have a Gaussian distribution by setting  $p$  and  $q$  to 2 and 10 respectively. In Figure 3 (b), the linear decrease shape of the beta-function is obtained by assigning a value of 0.01 to the parameter  $p$  and 1 to the parameter  $q$ . Figure 3 (c) also shows an exponential data distribution curve with  $p$  equal to 0.01 and  $q$  fixed to 1.

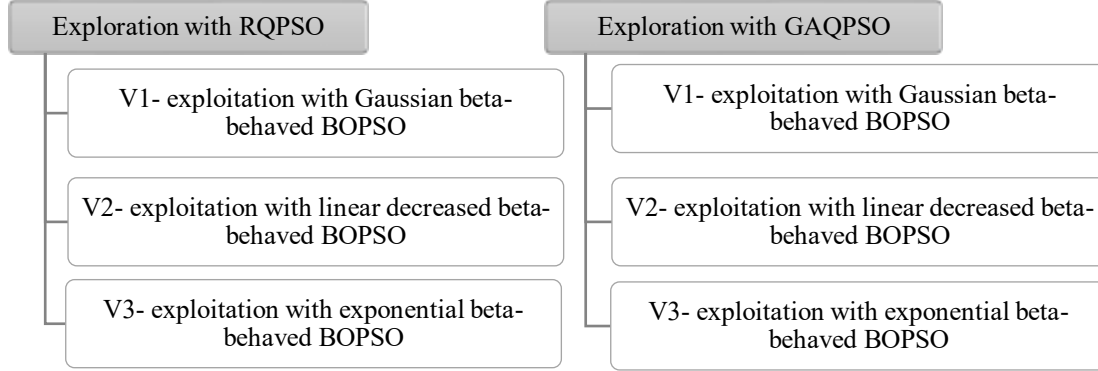


**Figure 3.** The Data Distribution Curve of Beta function with (a) Gaussian, (b) linear decrease, and (c) exponential distribution according to different configurations of  $p$  and  $q$ .

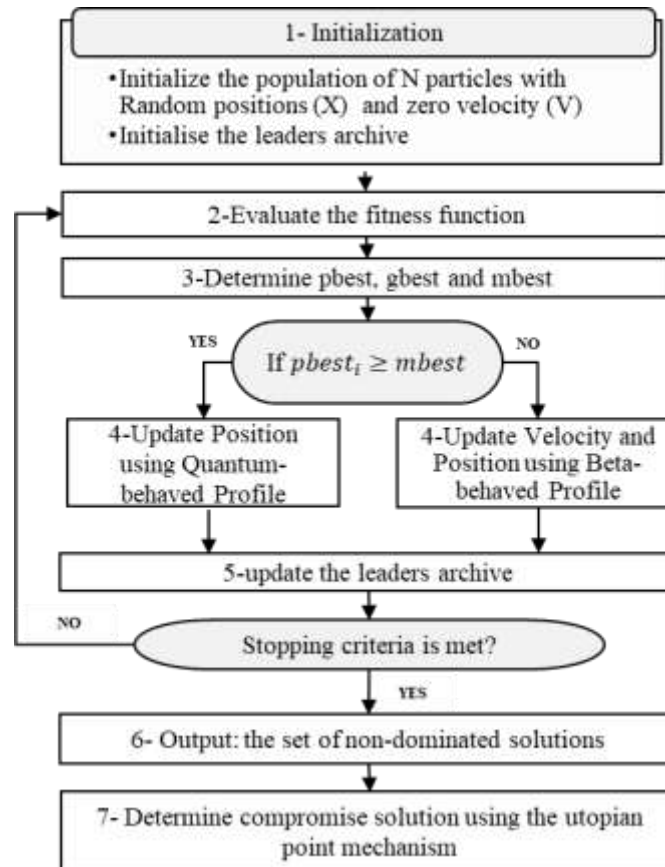
Taken into consideration the previous example that presents the diversification of the shapes that can be obtained by the beta function, six variants of QB-MOPSO approach are proposed based on GAQPSO and RQPSO for the exploration phase, and the three different distribution shapes of beta function for the exploitation enhancement. The six variants are developed to study the diversification of the data distribution of the new proposed QB-MOPSO algorithm, and all variants are illustrated in Figure 4 and detailed as follows:

- **In Revised Quantum Beta-behaved Multi-Objective Particle Swarm Optimization (RQB-MOPSO) system** the exploration profile is done using the update position of RQPSO algorithm as presented in Equation (8). However, three variants of beta exploitation profiles are done according to the parameter's configurations of the beta function in Equation (13). The three variants of RQB-MOPSO approach are as follows:
  - **RQB-MOPSO-V1:** exploration profile is done using RQPSO, and the exploitation profile is with a gaussian beta-behaved PSO.
  - **RQB-MOPSO-V2:** exploration profile is with RQPSO, and the exploitation phase is with linear decreased beta-behaved PSO.
  - **RQB-MOPSO-V3:** exploration phase is with RQPSO, and exploitation with exponential beta-behaved PSO.
- **In the Gaussian Quantum Beta-behaved Multi-Objective Particle Swarm Optimization (GAQB-MOPSO) system,** the particles positions are updated using Equation (9) of GAQPSO for the exploration step. According to the three-beta configurations in Equation (13) for the exploitation profile, the three variants of the GAQB-MOPSO system are as follows:
  - **GAQB-MOPSO-V1:** exploration is with GAQPSO, and exploitation is with Gaussian beta-behaved PSO.

- **GAQB-MOPSO-V2**: exploration is with GAQPSO, and exploitation is with a linear decreased beta function.
- **GAQB-MOPSO-V3**: exploration with GAQPSO, and exploitation with exponential beta-behaved PSO.



**Figure 4.** Six Variants of the Proposed QB-MOPSO Algorithm with (V1) Gaussian, (V2) linear decreased, and (V3) exponential beta function



**Figure 5.** The Flowchart of the Proposed Quantum Beta-behaved Multi-Objective Particle Swarm Optimization (QB-MOPSO) System.

Figure 5 illustrates the flowchart of the proposed QB-MOPSO approach. The details of the steps are as follows:

✓ **Step 1: initialization**

The first step aims to create a swarm of  $N$  particles with random positions  $X_i(t)$  and zero velocity  $V_i(t)$  vectors. Each particle  $p_i$ ,  $\forall i = 1 \dots N$  has defined in  $m$ -dimensional search space. The iterative optimization process is considered to evaluate the fitness function and to update the particle positions. At each iteration, all non-dominated solutions are stored in the leader's archive.

✓ **Step 2: fitness function evaluation**

At each iteration  $t$ , a predefined fitness function  $F(x, t)$  was evaluated.

✓ **Step 3: select pbest, gbest and mbest**

The global best solution (*gbest*) is selected randomly from the leader's archive. The personal best solution (*pbest*), is the best historical experience. The mean personal best solution (*mbest*) is determined using Equation (7).

✓ **Step 4: update the particles positions**

Compared to the existing PSO approaches, the proposed QB-MOPSO algorithm benefits from a new equation to update the particles positions. As presented in Equation (10), a new optimization equation is used with two optimization profiles, where the position  $X$  of each particle  $p_i$  is distributed symmetrically about the mean personal best position (*mbest*). The first optimization profile is for exploration phase using the Quantum-behaved equation of RQPSO and GAQPSO. However, the second optimization profile is based on Beta-behaved function for exploitation enhancement.

✓ **Step 5: Update the leader's archive**

At each iteration, all non-dominated solutions in the leader's archive ( $A$ ) are re-evaluated, and the dominated solutions are removed from the archive [48].

✓ **Step 6: Stopping criterion**

The QB-MOPSO system is stopped when the maximum number of iterations is met.

✓ **Step 7: Output of QB-MOPSO: determine the non-dominated solutions**

At the maximum number of iterations, a set of compromise solutions which are stored in the leader's archive ( $A$ ) are considered as the output of the proposal and denoted by Pareto Optimal Front (POF).

✓ **Step 8: Decision Making: determine the best compromise solution**

For decision making, the most known standard criterion in the optimization field is the use of the Utopian point mechanism [49], which is defined as an ideal infeasible solution that minimises the objective functions. After determining the utopian point, the Euclidian distance between this point and all non-dominated solutions in POF is computed. Then, the optimal particle with the smallest distance to the utopian point is selected as a compromise solution.

### 3.2 The Complexity Analysis of the Proposed QB-MOPSO Algorithm

Let us determine the complexity of the QB-MOPSO algorithm. The proposed QB-MOPSO algorithm aims to optimize a swarm of  $N$  particles, each particle is a candidate solution performed until a maximum number of iterations  $T_{max}$  is reached. First, at the iteration  $t = 1$  the initialization procedure is started including the following steps:



- Initialize random positions  $X_i$ ,  $\forall i = 1 \dots N$  with d-dimensional search space and velocities and takes  $O(N * d)$  times.
- The complexity time to evaluate the fitness function for  $N$  particles is equal to  $O(N * m)$ , where  $m$  is the number of objective functions.
- Apply the dominance operator to determine non-dominated solutions and stored in the leader's archive ( $A$ ) and takes  $O(N)$ .

Second, the main loop is executed until the maximum number of iterations  $T_{max}$  is reached. The running time of the QB-MOPSO algorithm consists of  $K$  iterative loops performing logarithmic statements and takes  $O(N * m \log(T))$  times. At each iteration, the above steps 2 to 6 are executed. The update of particle positions for the exploration or the exploitation profile is being preceded by determining the global best solution ( $gbest$ ) from the leader's archive ( $A$ ), the personal best position for each particle ( $pbest$ ), and the mean best particle ( $mbest$ ). Furthermore, the fitness function is evaluated during  $O(N * m)$ . At each time  $T$ , the leaders archive ( $A$ ) is updated and all dominated solutions are removed and takes  $O(N)$  times. Finally, the best compromise solution is determined using the Utopian point mechanism and the determination takes  $O(N)$  times. To sum up, the overall complexity of the proposed QB-MOPSO algorithm is equal to  $O(N * m \log(T))$ .

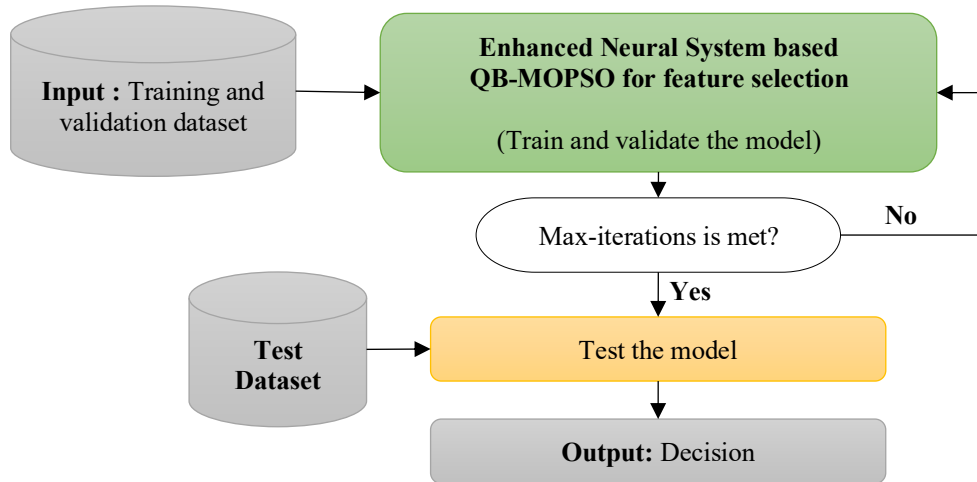
Based on previous work, Sun *et al.* [47] have concluded the high performance of Quantum behaved PSO (QPSO) compared to the traditional PSO algorithm [34], characterised with fewer control parameters and assumes a high level of convergence during the optimization process. However, the main advantage of the proposed QB-MOPSO algorithm is proved over their simplicity in terms of complexity which is equal to  $O(N * m \log(T))$ . The QB-MOPSO algorithm in this work benefits from two optimization profiles for exploration and exploitation phases. When, the dynamic switching profiles are the main properties of the proposed algorithm investigating a high flexibility to produce several types of data distributions. The quantum and beta behaved rules provide a higher level of convergence toward the global best solutions

### 3.3 The Neuro-QB-MOPSO Architecture for Fake Account Detection

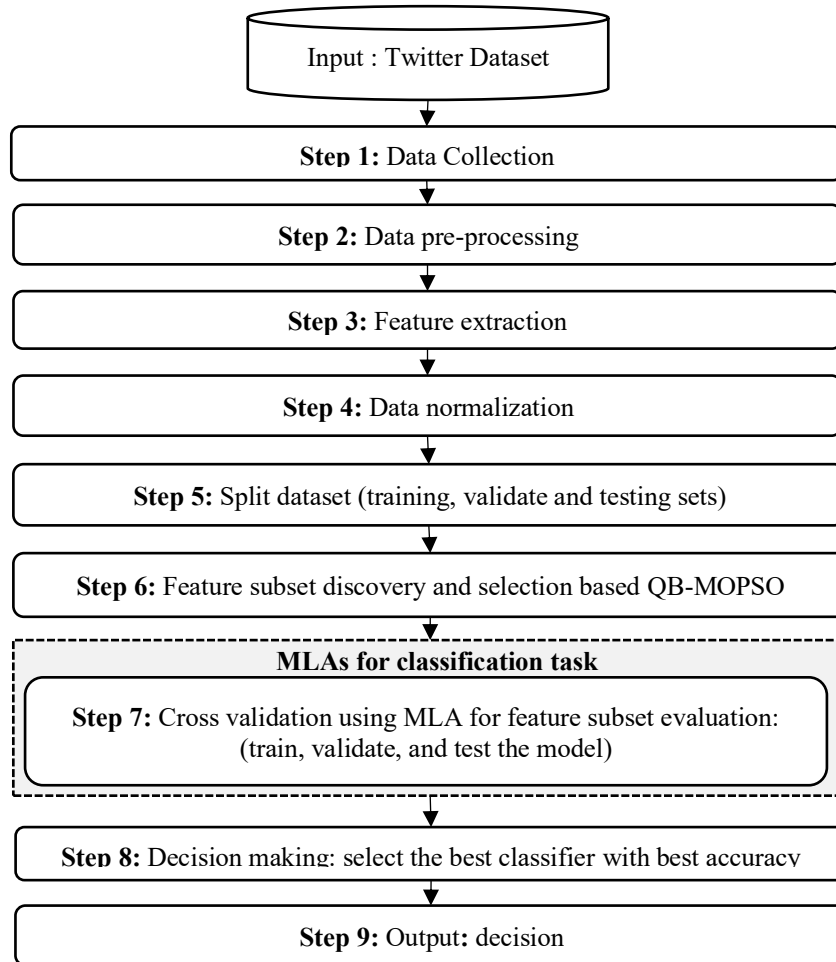
In this section presents the application of the proposed QB-MOPSO algorithm for Twitter fake accounts detection and denoted by the Neuro-QB-MOPSO Architecture.

#### a. General Description: Neuro-QB-MOPSO Architecture

The QB-MOPSO algorithm is proposed for pertinent features selection to detect fake accounts on Twitter. Figure 6 shows the Neuro-QB-MOPSO architecture. The Neuro-QB-MOPSO system takes a labelled Twitter dataset as input and performs the iterative process of the QB-MOPSO algorithm to determine the features to be selected. Each particle has a subset of selected features that are used in training and validating the machine learning model. At the maximum number of iterations, the model is tested with the best feature set that has the lowest error rate. Finally, the list of fake accounts is determined as the output of the proposed system. Four steps are involved in achieving the feature selection step, namely; dataset collection, data pre-processing, feature extraction, and feature analysis. Figure 7 illustrates the overall steps.



**Figure 6.** The Neuro-QB-MOPSO Architecture



**Figure 7.** The Overall Process of QB-MOPSO Algorithm for Deep Neural Architectures

Figure 7 shows the overall process for detecting fake accounts using the QB-MOPSO algorithm. The steps are as follows:

#### b. Dataset Collection

This step provides a starting point before diving into data exploration, and aims to select the main important columns from the original datasets. As shown in Table 2, two Twitter datasets proposed by Cresci *et al.* [31] are considered. A total number of 2910 Twitter accounts stored in both datasets. Dataset 1 contains 1982 Twitter accounts, while dataset 2 contains 928 accounts equally divided between human and social spam bot accounts.

**Table 2.** Cresci Datasets Properties

Properties	Dataset 1	Dataset 2	Total
Nb. accounts	1982	928	2910
Nb. tweets	4 061 598	2 628 181	6 689 779
Description	Genuine accounts + Social Spam Bot 1 (retweets of an Italian political candidate)	Genuine accounts + Social Spam Bot 3 (spammers of products on sale at Amazon.com)	-

#### c. Data Pre-processing

Several tasks are considered for text cleaning and presented as follows:

- Convert the corpus of tweets to lowercase.
- Removing the numbers from tweets using regular expressions to reduce the irrelevant features.
- Remove the set of symbols or punctuations.
- Remove white spaces from the tweet.
- Stop word removal: remove the common words in the language that do not carry a relevant meaning using natural language processing mechanisms.
- Stemming step to reduce the word to its stem forms using the Porter stemming algorithm that aims to remove the common morphological and inflexional endings from words (examples: users → user, profiling → profile).
- Lemmatization aims to reduce the word to the correct base forms using the lemmatization tool denoted by WordNet Lemmatize presented in Python Natural Language Toolkit (NLTK) library [50].

#### d. Features Extraction

In this study, the user profile properties and the content of tweets are the main information sources for feature extraction. As mentioned by Rostami and Karbasi [17], 46 original features are collected from Cresci datasets [10] using a set of standard statistical criteria such as entropy, and standard deviation. Table 3, has detailed 22 features extracted from the user profile information's. Table 4 presents 24 extracted features based on the tweets content.

**Table 3.** Extracted Features based on User profile properties

ID	Features based on User profile properties
F1	Follower count
F2	Follower count/Account Age
F3	Following count
F4	Following count/Account Age
F5	Follower count/ Following count
F6	Follower count/ Following count
F7	(2* Follower count) -Following count
F8	Follower count/ Follower + Following
F9	Favorite's count
F10	Favorite's count/ Account Age
F11	Tweet count
F12	Tweet count/Account Age
F13	List count
F14	List count/ Account Age
F15	Favorites count/Tweet count
F16	List count/ Follower count
F17	GEO Tag
F18	Retweet count
F19	Retweet count/Tweet count
F20	The consecutive Retweets interval mean
F21	The consecutive Retweets interval Standard deviation
F22	Number of times the tweets sent by the user are retweeted by other users

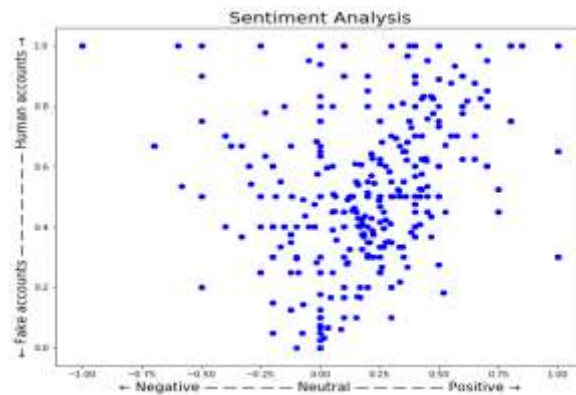
**Table 4.** Extracted Features based on Tweets Content

ID	Features based on tweets content
F23	Hashtag count
F24	Hashtag count/ Tweet count
F25	Hashtag-per-tweet Standard deviation
F26	Hashtag-per-tweet Entropy
F27	Tweets-with-Hashtags proportion
F28	The consecutive tweets interval means
F29	The consecutive tweets interval Standard deviation
F30	Link count
F31	Link count/ Tweet count
F32	Link-per-tweet standard deviation
F33	Link-per-tweet Entropy
F34	Tweets-with-Links proportion
F35	Mention count
F36	Mention count/Tweet count
F37	Mention-per-tweet Standard deviation
F38	Mention-per-tweet Entropy
F39	Tweets-with-Mentions proportion
F40	Reply count
F41	The consecutive Replies interval mean
F42	The consecutive Replies interval standard deviation
F43	Reply count/ Mention count
F44	The total number of the received likes
F45	Received likes count/ Tweet count
F46	Received like-per-tweet Standard deviation

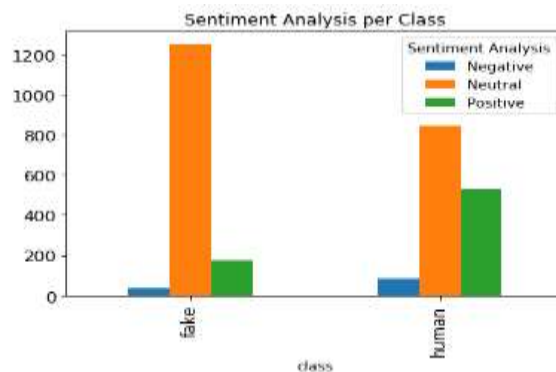
### e. Behavior Analysis of Fake and Human Web Users

The behavior analysis is considered to understand the attitude and the ethics of fake and human users on Twitter. Many features can be extracted from the tweets and denoted by content-based features extraction aiming to extract several features by parsing the content of each tweet such as the number of hashtags per tweet, the number of mentions per tweet, the length of the tweet and many others. In this sub-section, the sentiment analysis [51] process was first done. Sentiment analysis is an important topic in the field of Natural Language Processing (NLP) and presents a very high subject over many studies to detect negative, positive, and neutral sentiment presenting a subjective opinion based on text analysis. In this contribution, the corpus of tweets is used to produce the following labels (1: positive, 0: neutral and -1: negative).

Before starting the sentiment analysis task, the text pre-processing step is considered to clean the corpus of tweets and detailed in Step 2. Furthermore, the Text-Blob python library is used for tweet processing and sentiment analysis by computing two properties; polarity and subjectivity for each tweet. The two properties are presented as a float value in the range of  $[-1; 1]$  for the polarity property and  $[0; 1]$  for the subjectivity property. The two Cresci datasets [10] have 2910 online user accounts regrouped equally; 1455 for fake and 1455 for human accounts. Based on Figures 8 and 9, it is remarkable that the most important number of tweets are with a neutral sentiment. However, we can conclude that the human accounts can express their opinion and feelings in the corpus of the tweet compared with fake users which have a large number of neutral opinions.



**Figure 8:** Sentiment analysis using polarity and subjectivity for fake and human accounts.

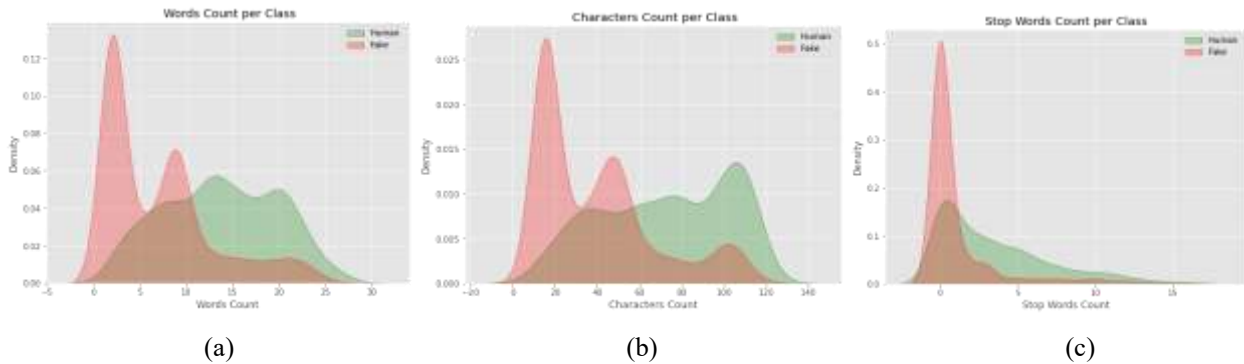


**Figure 9.** Classification of sentiment analysis for human and fake profiles.

Table 5 presents the number (Nb.) and the proportion in percentage of the users accounts according to the results of the sentiment analysis process. For fake accounts, there are 2.41%, 85.84% and 11.75% for negative, neutral, and positive opinions respectively compared to the human users which have 5.49%, 58.08% and 36.43%.

**Table 5.** Comparative results of sentiment analysis for bot and human accounts on Twitter.

Sentiment	Fake		Human	
	Nb.	%	Nb.	%
Negative	35	<b>2.41</b>	80	<b>5.49</b>
Neutral	1249	<b>85.84</b>	845	<b>58.08</b>
Positive	171	<b>11.75</b>	530	<b>36.43</b>
Total	1455	100	1455	100



**Figure 10.** (a) Words Count, (b) Characters Count, and (c) Stop Words Count of Fake and Human Tweets Content.

The plots in Figure 10 illustrates the words count, the number of characters, and the number of stop words per class using the Tweet content. It can be remarkable that most fake users have the habit of writing a Tweet corps of 5 to 10 words composed of 10 to 50 characters. However, the tweet corpus of most human users is composed of 15 to 25 words and 80 to 110 characters. A spam-bot or fake user does not have the habit of using a lot of stop words compared with human user.



**Figure 11.** Words Cloud of (a) Fake and (b) Human Tweets Content.

Figure 11 presents the words cloud to visualize the representation of the word frequency of fake and human tweets content, and highlight the words which are frequently used.

#### f. Feature selection based QB-MOPSO algorithm

In this sub-section, the Neuro-QB-MOPSO system is applied to the feature selection problem to minimize simultaneously the features dimensionality using Equation (14) and the classification error rate using Equation (15). The main steps of the Neuro-QB-MOPSO system are detailed as follows:

In the QB-MOPSO algorithm, the position of each particle is n-dimensional equal to the number of features in the dataset. In addition, the position Binary Encoding ( $XBE(t)$ ) is added and initialized to present the bits of selected features that monitor the evaluation of the classifier.

##### ✓ Feature subset discovery

At each iteration ( $t$ ), the feature subset discovery step generates the candidate feature subset. It is a search procedure started with all available features in the dataset. The search step aims to find the best subset of features investigated to evaluate the fitness function of each particle.

##### ✓ Feature subset selection

The binary position ( $XBE_i(t)$ ) of each particle  $p_i$  is computed using the Sigmoid Transfer Function (STF) in Equation (16).

$$STF_i(t) = \frac{1}{1+e^{-X_i(t)}} \quad (16)$$

where:  $X_i(t)$  is the position of the particle  $i$ .

The binary bits ( $XBE_i(t)$ ) is determined according to Equation (17):

$$XBE_i(t) = \begin{cases} 1 & \text{if } STF_i(t) \geq rand \\ 0 & \text{otherwise} \end{cases} \quad (17)$$

In the vector of binary position ( $XBE_i(t)$ ), “1” represents the feature is selected and “0” otherwise.

##### ✓ Subset evaluation

For each particle, the subset of selected features is evaluated using the two objective functions provided in Equations (14) and (15) to train the classification model. This step is optimized until a maximum number of iterations. At each iteration, all the non-dominated solutions are stored in the leader’s archive ( $A$ ).

##### ✓ Decision making for best features subset selection

At the maximum number of iterations, the Pareto Optimal Front (POF) is used to determine the Utopian Point and the compromise solution presenting the best set of multiple features. Finally, the best subset of selected features is considered to test and determine the performance of MLAs.

##### ✓ Stopping criterion

The optimization process of QB-MOPSO system is stopped when the maximum number of iterations is met.

##### ✓ The output of QB-MOPSO for fake account detection

Determine the compromise solution, the subset of the selected features and best classifier with best accuracy.

As shown above, the application of QB-MOPSO to detect fake accounts involves nine steps, which are illustrated in Figure 7: Data collection, data pre-processing, feature extraction, data normalisation and splitting the dataset into training, validation, and testing sets. Then, the QB-MOPSO algorithm is used to determine the best feature set when all particles are considered for training and validating the machine learning model. However, the best feature set

selected by the compromise solution is considered for testing the model and has the highest accuracy rate. Finally, the classifier with the best accuracy is selected for decision making and the list of fake accounts is generated.

#### 4. Experimental Study

This section presents the experimental study. Subsection 6.1 outline the state-of-the-art methods and explains the performance metrics used. Subsection 6.2 describes the parameter settings. The quantitative results and discussion are presented in subsection 6.3.

##### 6.1 Preliminary and Performance Metrics

In this experimental study, all variants of the QB-MOPSO approach are implemented on a personal computer with 8 GB of RAM, 1 To and i7 Intel processors using MATLAB programming platform. The proposed QB-MOPSO algorithm is compared with the state-of-the-art methods proposed by the following authors: Rostami and Karbasi [17], Ahmed and Abulaish [19], Davis *et al.* [21], Cresci *et al.* [22], Yang *et al.* [23], Miller *et al.* [24] using the two Twitter datasets provided by Cresci *et al.* [10] detailed in Table 2. The experimental results of this study are compared to the methods discussed by Rostami *et al.* [17] using three machine learning algorithms (MLAs) namely: Random Forest (RF), Naïve Bayes (NB), and Support Vector Machine (SVM). 10-fold cross-validation technique is used to compute the performance criteria of RF, NB, and SVM algorithms. However, hidden layers are considered for the feedforward Neural Network (NN) algorithm, and the split of the datasets is as follows; 70% for training, 15% for validation and 15% for testing. For each MLA, we aim to compute the confusion matrix for performance measurement. In this case of study, three performance criteria are considered:

- ✓ **Accuracy (Acc.)** or classification rate presents the percent of the correct classified samples computed using Equation (18).

$$Accuracy = \frac{TP+TN}{TN+FN+TP+FP} \quad (18)$$

where, TP, TN, FN, and FP have been easily determined from the confusion matrix in Table 6.

**Table 6.** the confusion matrix

True class	Predicted class		
		True	False
	True	True Positive (TP)	False Positive (FN)
	False	True Negative (TN)	False Negative (FP)

- ✓ **F-Measure** is computed using Equation (19) based on the precision and recall criteria.

$$F - \text{measure} = \frac{2 \times \text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}} \quad (19)$$

where; precision and recall are respectively equal to  $TP / (TP + FP)$  and  $TP / (TP + FN)$ .



- ✓ **Matthew's correlation coefficient (MCC)** is the most important criteria for classification performance, and calculated using Equation (20).

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (20)$$

## 6.2 Orthogonal Experimental Design for the QB-MOPSO Algorithm

The parameter design of six variants of the proposed QB-MOPSO algorithm is done based on the Taguchi method [52]. The QB-MOPSO algorithm has four parameters (control factors) which are tested for sensitivity analysis, each parameter has two-level factors and listed as follows:

- ✓ The population size is fixed to 50 and 100.
- ✓ The maximum number of iterations is fixed to 50 and 100.
- ✓ The parameter  $p$  of the Beta function is fixed to 0.01 and 2.
- ✓ The parameter  $q$  of the Beta function is fixed to 1 and 10.
- ✓ Parameters of both Revised Quantum PSO (RQPSO) and Gaussian QPSO (GAQPSO) algorithms are considered as the original publications in [11] and [12] and detailed in Equations (8) and (9).

As mentioned in Table 7, the application of the Taguchi orthogonal arrays has identified an array of  $L_8(2^4)$  with only 8 best runs for different combination of parameters design with two-level design of 4 control factors.

**Table 7.** Taguchi Array Design  $L_8(2^4)$  of QB-MOPSO Variants

ID Run	Population size	Max-iterations	Parameters of Beta function	
			$p$	$q$
1	50	50	0,01	1
2	50	50	2	10
3	50	100	0.01	10
4	50	100	2	1
5	100	50	0.01	10
6	100	50	2	1
7	100	50	0.01	1
8	100	100	2	10

Based on the Taguchi array design, the six QB-MOPSO variants are tested for 8 independent runs and the best obtained results of the GAQB-MOPSO algorithm are reported in Table 8 and for the RQB-MOPSO algorithm are reported in Table 9 presenting the accuracy of four machine learning algorithms namely; RF, SVM, NN and NB using the datasets 1 and 2.

**Table 8.** The Accuracy Rate of MLAs using GAQB-MOPSO algorithm using Taguchi Method for Datasets 1 and 2.

MLAs	ID Run	Population size	Max-iterations	Parameters of Beta function		Accuracy %	
				p	q	Dataset 1	Dataset 2
NB	1	50	50	0,01	1	94.00	85.80
	2	<b>50</b>	<b>50</b>	<b>2</b>	<b>10</b>	<b>95.20</b>	85.90
	3	50	100	0,01	10	94.90	85.10
	4	<b>50</b>	<b>100</b>	<b>2</b>	<b>1</b>	94.10	<b>86.00</b>
	5	100	50	0,01	10	93.60	85.10
	6	100	50	2	1	94.00	62.00
	7	100	50	0,01	1	95.10	85.10
	8	100	100	2	10	93.60	62.00
NN	1	<b>50</b>	<b>50</b>	<b>0,01</b>	<b>1</b>	99.10	<b>93.50</b>
	2	50	50	2	10	98.90	85.20
	3	50	100	0,01	10	98.90	89.70
	4	50	100	2	1	98.90	84.10
	5	100	50	0,01	10	98.60	89.30
	6	100	50	2	1	98.20	90.60
	7	<b>100</b>	<b>50</b>	<b>0,01</b>	<b>1</b>	<b>99.19</b>	91.10
	8	<b>100</b>	<b>100</b>	<b>2</b>	<b>10</b>	<b>98.50</b>	96.44
RF	1	50	50	0,01	1	97.40	96.33
	2	<b>50</b>	<b>50</b>	<b>2</b>	<b>10</b>	<b>97.43</b>	96.98
	3	<b>50</b>	<b>100</b>	<b>0,01</b>	<b>10</b>	97.41	<b>97.62</b>
	4	50	100	2	1	97.19	96.98
	5	100	50	0,01	10	96.76	97.52
	6	100	50	2	1	97.19	97.09
	7	100	50	0,01	1	96.98	96.65
	8	100	100	2	10	96.33	93.30
SVM	1	50	50	0,01	1	94.50	93.30
	2	<b>50</b>	<b>50</b>	<b>2</b>	<b>10</b>	92.90	<b>94.80</b>
	3	50	100	0,01	10	92.90	92.80
	4	50	100	2	1	95.40	93.00
	5	<b>100</b>	<b>50</b>	<b>0,01</b>	<b>10</b>	<b>95.50</b>	92.90
	6	<b>100</b>	<b>50</b>	<b>2</b>	<b>1</b>	93.10	<b>94.80</b>
	7	100	50	0,01	1	92.90	92.90
	8	100	100	2	10	92.90	93.20

**Table 9.** The Accuracy Rate of MLAs using RQB-MOPSO algorithm using Taguchi Method for Datasets 1 and 2.

MLAs	ID Run	Population size	Max-iterations	Parameters of Beta function		Accuracy %	
				p	q	Dataset 1	Dataset 2
NB	1	50	50	0.01	1	94.00	85.00
	2	50	50	2	10	94.00	84.10
	3	50	100	0.01	10	93.60	84.10
	4	50	100	2	1	94.10	85.80
	5	100	50	0.01	10	94.00	83.20
	6	<b>100</b>	<b>50</b>	<b>2</b>	<b>1</b>	94.00	<b>86.00</b>
	7	<b>100</b>	<b>50</b>	<b>0.01</b>	<b>1</b>	<b>95.00</b>	85.80
	8	100	100	2	10	93.70	83.20
NN	1	50	50	0.01	1	98.10	90.5
	2	<b>50</b>	<b>50</b>	<b>2</b>	<b>10</b>	<b>98.90</b>	90.00
	3	50	100	0.01	10	98.50	91.60
	4	50	100	2	1	98.30	83.30
	5	<b>100</b>	<b>50</b>	<b>0.01</b>	<b>10</b>	98.50	<b>94.90</b>
	6	100	50	2	1	98.80	91.60
	7	<b>100</b>	<b>50</b>	<b>0.01</b>	<b>1</b>	<b>98.90</b>	90.50
	8	100	100	2	10	98.50	92.20
RF	1	<b>50</b>	<b>50</b>	<b>0.01</b>	<b>1</b>	97.09	<b>97.30</b>
	2	50	50	2	10	96.76	96.01
	3	50	100	0.01	10	96.87	96.98
	4	50	100	2	1	95.90	96.98
	5	<b>100</b>	<b>50</b>	<b>0.01</b>	<b>10</b>	97.62	<b>97.30</b>
	6	100	50	2	1	97.30	96.33
	7	100	50	0.01	1	96.98	96.98
	8	100	100	2	10	97.09	96.22
SVM	1	50	50	0.01	1	92.90	92.90
	2	<b>50</b>	<b>50</b>	<b>2</b>	<b>10</b>	92.60	<b>93.30</b>
	3	50	100	0.01	10	92.80	93.20
	4	<b>50</b>	<b>100</b>	<b>2</b>	<b>1</b>	<b>93.20</b>	93.00
	5	100	50	0.01	10	92.80	92.70
	6	100	50	2	1	93.00	93.20
	7	100	50	0.01	1	93.10	92.80
	8	<b>100</b>	<b>100</b>	<b>2</b>	<b>10</b>	<b>93.20</b>	92.90

### 6.3 Parameters Settings

Based on the orthogonal experimental design using the Taguchi method, the best parameters setting of all variants of the QB-MOPSO approach are resumed in Table 10 presenting the best configuration to obtain the best accuracy rate for all MLAs using datasets 1 and 2.

**Table 10.** Parameters Setting for RQB-MOPSO and RQB-MOPSO variants with (V1) Gaussian, (V2) linear decreased, and (V3) exponential beta function.

Profiles of QB-MOPSO	Parameters	RQB-MOPSO			GAQB-MOPSO		
		V1	V2	V3	V1	V2	V3
Quantum-behaved PSO	$\varphi_1$ $\varphi_2$ u	<i>rand</i> (0,1)			<i>N</i> (0,1)		
Beta-behaved PSO	p	2	0.01	0.01	2	0.01	0.01
	q	10	1	10	10	1	10
Swarm size	100						
Max-number of iterations	50						

### 6.4 Results and Discussion

This section presents experimental results with all the available features, as well as comparative results using a subset of selected features based on QB-MOPSO variants.

#### ✓ Quantitative results using all features

Four MLAs are first tested using all the available features in both datasets. The purpose of this study is to illustrate the importance of using a small subset of features and their impact on the classification accuracy. In Table 11, the quantitative results are shown using 46 original features. I can be evidently seen that NN is the best classifier, with an accuracy rate of 98.89% using the first dataset. Nevertheless, the random forest is the best classifier for the second dataset with an accuracy rate of 96.12%.

**Table 11.** Performance criterions (%) of MLAs using all features.

	MLAs	Performance criterions (%) with all features		
		Acc.	F-Measure	MCC
Dataset 1	Random Forest	97.78	97.78	95.56
	Naïve Bayes	84.66	83.37	70.18
	SVM	97.83	97.87	95.74
	Neural Network	<b>98.89</b>	<b>98.89</b>	<b>97.78</b>
Dataset 2	Random Forest	<b>96.12</b>	<b>96.12</b>	<b>92.24</b>
	Naïve Bayes	84.91	84.81	69.83
	SVM	93.32	93.11	86.79
	Neural Network	94.82	94.71	89.73

#### ✓ Quantitative results of QB-MOPSO variants using a subset of selected features on datasets 1 and 2

The six variants of the QB-MOPSO system are tested using both datasets 1 and 2. Based on data distribution types for exploration and exploitation profiles, there are three RQB-MOPSO variants (V1, V2 and V3) and three GAQB-MOPSO variants (V1, V2 and V3). In all variants, particles “fly” symmetrically to the center-of-gravity, which is the

mean *pbest* presenting the global attractor particle. Particles whose positions are greater or equal to the mean solution position are considered in the exploration profile, otherwise they are considered in the exploitation profile. Within an exploration profile, a RQB-MOPSO algorithm updates particle positions according to a uniform random distribution between 0 and 1, but a GAQB-MOPSO algorithm updates them according to a gaussian distribution with zero mean and unit variance. In addition, the particles positions in the exploitation profile are updated by using beta functions for gaussian (V1), linear decrease (V2), and exponential (V3) distributions. The classification performance of NB, RF, SVM, and NN classifiers is detailed in Table 12, along with the dimensionality subset of selected features determined by the different variants of the QB-MOPSO system. It is remarkable that QB-MOPSO variants are more competitive than the state-of-the-art methods using dataset 1.

The reported performance criteria of the NN classifier have demonstrated that GAQB-MOPSO (V2) is the best approach for fake account detection using dataset 1 with the highest accuracy rate of 99.19% using 32 selected features. Moreover, the comparative results in Table 13 have shown the superiority of GAQB-MOPSO (V2) using dataset 2 with an accuracy rate of 97.52 % with 25 selected features. The first 32 features are divided equally as follows; 16 features are selected based on the user profile properties which are presented in Table 3 (ID: F2, F4, F5, F6, F7, F9, F12, F14, F15, F16, F17, F18, F19, F20, F21, F22) presenting 73% of all user profile properties, and 16 features are selected based on tweet content as presented in Table 4 (ID: F27, F28, F29, F30, F31, F32, F35, F36, F38, F40, F41, F42, F43, F44, F45, F46) presenting 67% of all features based on the tweet content. For the second 25 selected features, 10 features based on the profile's properties (45%) are selected and their ID in Table 3 are as follows (ID: F1, F2, F4, F5, F7, F9, F15, F16, F18, F21), and 15 features are based on the tweet content from Table 4 (ID: F23, F24, F26, F27, F29, F31, F32, F34, F35, F36, F37, F38, F43, F44, F46) presenting 63% of the features-based content.

#### ✓ Comparative Results of the best variant “GAQB-MOPSO (V2)” versus State-of-the-art Methods

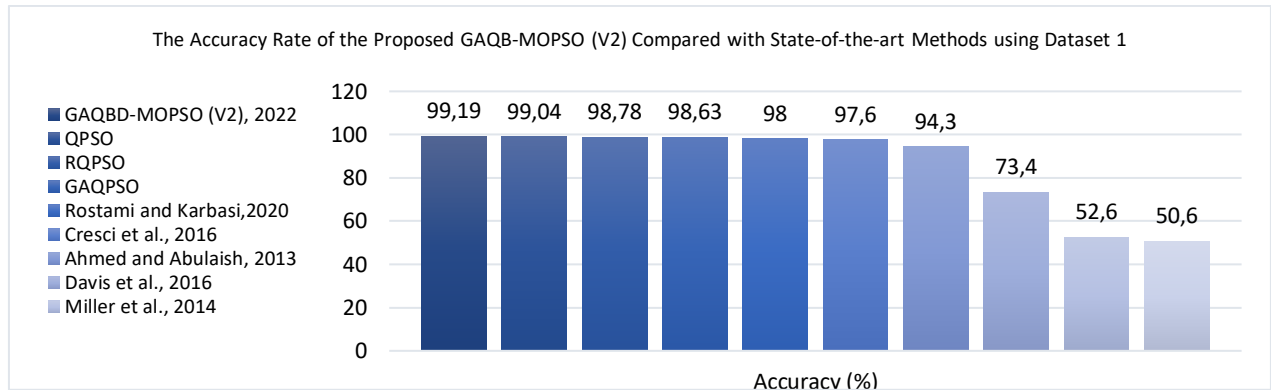
Tables 12 and 13 provide the classification performance of the state-of-the-art methods compared with all variants of the proposed QB-MOPSO algorithm. For both datasets, QB-MOPSO can assume a high accuracy rate compared with the supervised methods proposed by Rostami and Karbasi [17], Davis *et al.* [21] and Yang *et al.* [23] as well as the unsupervised approaches proposed by Ahmed and Abulaish [19], Cresci *et al.* [22], and Miller *et al.* [24]. The proposed QB-MOPSO presents high accuracy when using a supervised NN classifier, and it can identify 32 pertinent features from 46 original features when using dataset 1. Furthermore, only 25 features are selected with the GAQB-MOPSO (V2) system in RF classifier based on the dataset 2. In [17] Rostami and Karbasi, used the Minimum Redundancy –Maximum Relevance algorithm (mRMR) as feature selection technique, and had the ability to select 8 and 7 pertinent features using dataset 1 and dataset 2 respectively.

The selected features are based on tweets content, and SVM classifier is the top performer with an MCC equal to 96.06% for dataset 1, and 94.19% for dataset 2. Despite the minimum number of selected features (8 and 7) in [17], QB-MOPSO can achieve a high MCC (98.39% and 95.06%) with 32 and 25 optimal features using both datasets for the NN and RF respectively. In [22], Cresci *et al.* proposed a DNA fingerprinting method and achieve an MCC 95.20% and 86.70% regarding 14 features [19]. Davis *et al.* [21], proposed the BotOrNot system and achieved a classification rate of 17.4% on dataset 1 and 37.8% on dataset 2, using a random forest classifier with 1000 Twitter feature account.

In [24], Miller et al. have considered 126 features to test their model, and does not assume a good result compared to all methods. Also, Ahmed and Abulaish [19] have proposed Graph clustering and Community Detection methods and 14 generic statistical features are selected to test the unsupervised model. Finally, Yang et al. [23], have proposed an empirical analysis of profile-based feature evasion tactics and content-based feature evasion tactics using 25 features. However, it fails to obtain a good classification accuracy (MCC=4.3% with dataset 1 and MCC=28.7% with dataset 2). Compared with all methods, QB-MOPSO can achieve a good performance using a dynamic feature selection according to quantum weights and the diversification of beta profiles, which are encoded using the sigmoid function.

**Table 12.** Comparative Results of **QB-MOPSO** variants with (V1) Gaussian, (V2) linear decreased, and (V3) exponential beta function compared with state-of-the-art methods using **dataset 1**.

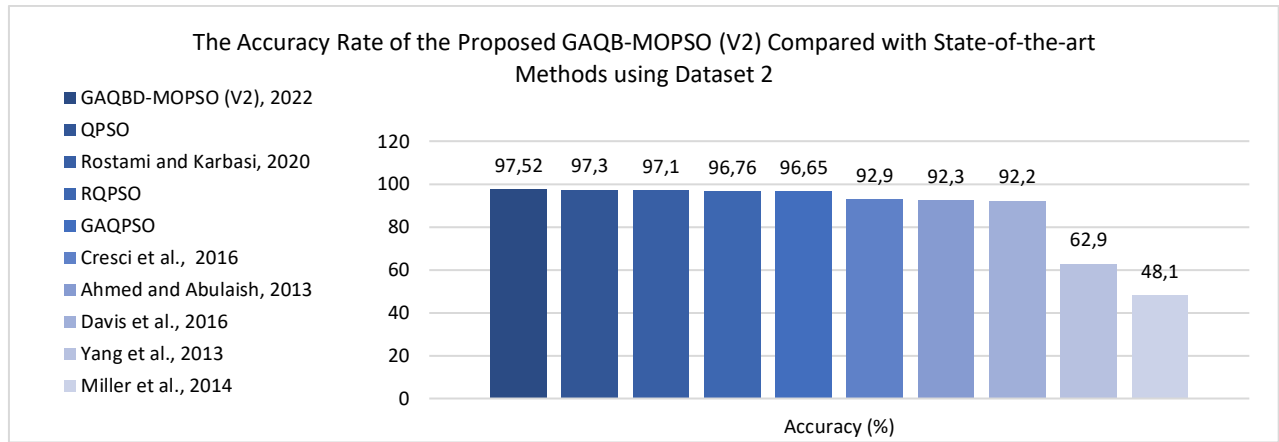
	Compared Approaches		Selected features	Performance criterions (%) using the best features subset									
				TP	TN	FP	FN	Precision	Recall	Specificity	Acc.	F-Measure	MCC
NB	RQB-MOPSO	V1	14	812	929	179	62	81.94	92.91	83.84	87.41	87.08	76.21
		V2	30	949	936	42	55	95.76	94.52	95.71	95.11	95.14	90.22
		V3	7	790	884	201	107	79.71	88.07	81.47	84.46	83.69	69.23
	GAQB-MOPSO	V1	31	950	936	41	55	95.86	94.52	95.80	95.15	95.19	90.32
		V2	21	931	934	60	57	93.94	94.23	93.96	94.09	94.08	88.19
		V3	19	922	935	69	56	93.03	94.27	93.12	93.69	93.65	87.39
SVM	RQB-MOPSO	V1	32	967	990	24	1	97.57	99.89	97.63	98.73	98.72	97.50
		V2	26	969	991	22	0	97.78	100	97.82	98.89	98.87	97.80
		V3	35	966	990	25	1	97.47	99.89	97.53	98.68	98.67	97.40
	GAQB-MOPSO	V1	26	970	990	21	1	97.88	99.89	97.92	98.89	98.87	97.79
		V2	25	967	991	24	0	97.57	100	97.63	98.78	98.77	97.60
		V3	30	968	990	23	1	99.89	99.89	95.83	98.78	99.89	97.60
RF	RQB-MOPSO	V1	14	975	985	16	6	98.38	99.38	98.40	98.89	98.88	97.78
		V2	12	970	991	21	0	97.88	100	97.92	98.94	98.92	97.90
		V3	13	975	986	16	5	98.38	99.48	98.40	98.94	98.93	97.88
	GAQB-MOPSO	V1	23	980	978	11	13	98.89	98.69	98.88	98.78	98.79	97.57
		V2	18	975	990	16	1	98.38	99.89	98.40	99.14	99.13	98.29
		V3	17	974	991	17	0	98.28	100	98.31	99.14	99.13	98.29
NN	RQB-MOPSO	V1	25	972	991	19	0	98.08	100	98.11	99.04	99.03	98.10
		V2	24	964	991	27	0	97.27	100	97.34	98.63	98.61	97.31
		V3	18	955	990	36	1	96.36	99.89	96.49	98.13	98.09	96.32
	GAQB-MOPSO	V1	27	959	990	32	1	96.77	99.89	96.86	98.33	98.30	96.71
		V2	32	976	990	15	1	98.48	99.89	98.50	99.19	99.18	98.39
		V3	17	974	991	17	0	98.28	100	98.31	99.14	99.13	98.29
State-of-the art methods	Rostami and Karbasi [17], 2020		-	-	-	-	-	98.00	98.10	98.00	98.00	98.00	96.06
	Cresci <i>et al.</i> [22], 2016		-	-	-	-	-	98.20	97.20	98.10	97.60	97.70	95.20
	Davis <i>et al.</i> [21], 2016		>1000	-	-	-	-	47.10	20.80	91.80	73.40	28.80	17.40
	Miller <i>et al.</i> [24], 2014		126	-	-	-	-	55.50	35.80	69.80	52.60	43.50	5.90
	Ahmed and Abulaish [19], 2013		-	-	-	-	-	94.50	94.40	94.50	94.30	94.40	88.60
	Yang <i>et al.</i> [23], 2013		25	-	-	-	-	56.30	17.00	86.00	50.60	26.10	4.30



**Figure 12.** The Accuracy Rate of the GAQB-MOPSO (V2) Compared with state-of-the-art Methods using Dataset 1

**Table 13.** Comparative Results of QB-MOPSO variants with (V1) Gaussian, (V2) linear decreased, and (V3) exponential beta function compared with state-of-the-art methods using **dataset 2**.

	Compared Approaches		Selected features	Performance criterions (%) using the best features subset									
				TP	TN	FP	FN	Precision	Recall	Specificity	Acc.	F-Measure	MCC
NB	RQB-MOPSO	V1	19	443	347	21	117	95.47	79.11	94.29	85.13	86.52	71.81
		V2	9	436	353	28	111	93.97	79.71	92.65	85.02	86.25	71.19
		V3	17	444	354	20	110	95.69	80.14	94.65	85.99	87.22	73.37
	GAQB-MOPSO	V1	15	446	351	18	113	96.12	79.78	95.12	85.88	87.19	73.32
		V2	29	453	149	11	315	97.62	58.98	93.12	64.87	73.53	39.36
		V3	21	444	336	20	128	95.68	77.62	94.38	84.05	85.71	70.02
SVM	RQB-MOPSO	V1	8	437	404	27	60	94.18	87.92	93.73	90.62	90.94	81.45
		V2	26	445	421	19	43	95.90	91.18	95.68	93.31	93.48	86.75
		V3	24	447	418	17	46	96.33	90.66	96.09	93.21	93.41	86.59
	GAQB-MOPSO	V1	21	445	417	19	47	95.90	90.44	95.64	92.88	93.09	85.93
		V2	20	444	419	20	45	95.68	90.79	95.44	92.99	93.17	86.11
		V3	24	445	421	19	43	95.90	91.18	95.68	93.31	93.48	86.75
RF	RQB-MOPSO	V1	20	443	460	21	4	95.47	99.10	95.63	97.30	97.25	94.67
		V2	22	449	451	15	13	96.76	97.18	96.78	96.98	96.97	93.96
		V3	14	455	448	9	16	98.06	96.60	98.03	97.30	97.32	94.62
	GAQB-MOPSO	V1	19	430	462	34	2	92.67	99.53	93.14	96.12	95.98	92.46
		V2	25	445	460	19	4	95.90	99.10	96.03	97.52	97.48	95.09
		V3	33	441	460	23	4	95.04	99.10	95.23	97.09	97.02	94.26
NN	RQB-MOPSO	V1	16	440	384	24	80	94.82	84.61	94.11	88.79	89.43	78.15
		V2	16	443	409	21	55	95.47	88.95	95.11	91.81	92.09	83.84
		V3	21	453	383	11	81	97.62	84.83	97.20	90.08	90.78	81.10
	GAQB-MOPSO	V1	27	445	420	19	44	95.90	91.00	95.67	93.21	93.38	86.54
		V2	22	443	402	21	62	95.47	87.72	95.03	91.05	91.43	82.43
		V3	31	454	438	26	10	94.58	97.84	94.39	96.12	96.18	92.29
State-of-the art methods	Rostami and Karbasi [17], 2020		-	-	-	-	-	96.50	97.90	96.40	97.10	97.10	94.19
	Cresci <i>et al.</i> [22], 2016		-	-	-	-	-	100	85.80	100	92.90	92.30	86.70
	Davis <i>et al.</i> [21], 2016		>1000	-	-	-	-	63.50	95.00	98.10	92.20	76.10	73.80
	Miller <i>et al.</i> [24], 2014		126	-	-	-	-	46.70	30.60	65.40	48.10	37.00	-4.30
	Ahmed and Abulaish [19], 2013		-	-	-	-	-	91.30	93.50	91.20	92.30	92.30	84.70
	Yang <i>et al.</i> [23], 2013		25	-	-	-	-	72.70	40.90	84.80	62.90	52.40	28.70



**Figure 13.** The Accuracy Rate of GAQB-MOPSO (V2) Compared with state-of-the-art Methods using Dataset 2.

Table 14 shown a comparative result of the proposed GAQB-MOPSO (V2) versus the three standard approaches to quantum-behaved PSO namely; QPSO, RQPSO and GAQPSO. To study the performance of the novel proposal compared with other original systems (QPSO, RQPSO and GAQPSO). All algorithms are executed in the same condition with the NN for dataset 1 and the RF for dataset 2. Of course, it is remarkable that GAQB-MOPSO (V2) is the winner and has achieved the highest accuracy rates with NN classifier using dataset 1 and RF using dataset 2. Figures 13 and 14 have shown that GAQB-MOPSO (V2) has the highest accuracy rate for detecting fake accounts compared to the existing methods.

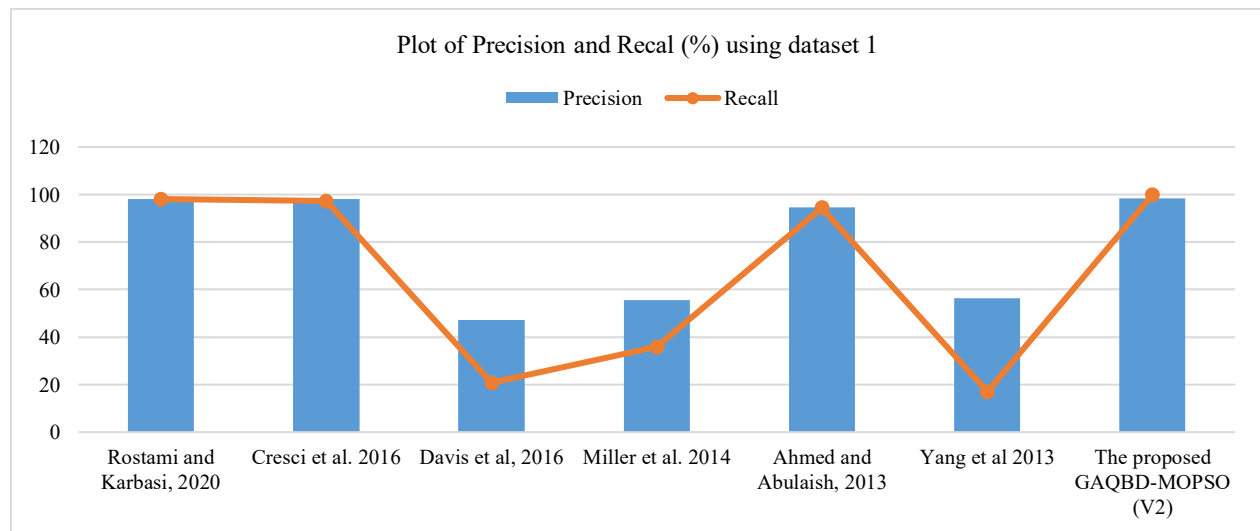
**Table 14.** Performance criterions (%) of the best classifiers (Random Forest and Neural Network) for the proposed GAQB-MOPSO (V2) versus the standard QPSO, RQPSO, and GAQPSO.

Datasets	Compared Approaches	Selected features	Performance criterions (%) using the best features subset									
			TP	TN	FP	FN	Precision	Recall	Specificity	Acc.	F-Measure	MCC
Dataset 1	GAQB-MOPSO-V2 (NN)	32	976	990	15	1	98.48	99.89	98.50	<b>99.19</b>	99.18	<b>98.39</b>
	QPSO (NN)	33	974	989	17	2	98.28	99.79	98.31	99.04	99.03	98.09
	RQPSO (NN)	22	967	991	24	0	97.57	100	97.63	98.78	98.77	97.60
	GAQPSO (NN)	20	966	989	25	2	97.47	99.79	97.53	98.63	98.62	97.30
Dataset 2	GAQB-MOPSO-V2 (RF)	25	445	460	19	4	95.90	99.10	96.03	<b>97.52</b>	97.48	<b>95.09</b>
	QPSO (RF)	25	443	460	21	4	95.47	99.10	95.63	97.30	97.25	94.67
	RQPSO (RF)	19	448	450	16	14	96.55	96.96	95.56	96.76	96.76	93.53
	GAQPSO (RF)	29	437	460	27	4	94.18	99.09	94.45	96.65	96.57	93.43

The state-of-the-art methods compared with the proposed QB-MOPSO algorithm have been divided into supervised and unsupervised methods for classifying spambots. In supervised methods, Yang *et al.* [23] proposed a spambot detection system based on machine learning algorithms to predict human and spambot accounts. In addition, Davis *et al.* [21] proposed the BotOrNot Blackbox platform. The analysis of the results in Figure 15 shows that both the methods in [21] and [23] failed in classification and most bot accounts were classified as human with a recall rate less than 50. For dataset 2, Figure 16 shows that the precision and recall values of the system proposed by Davis *et al.* [21] are unbalanced, resulting in lower accuracy of the classification model, and that the system of Yang *et al.* [23] fails with a recall rate below 50. In unsupervised methods, Miller *et al.* [24] proposed a stream clustering model based on DenStream [53] and StreamKM++ [54] clustering algorithms for spambot detection to determine the cluster of the

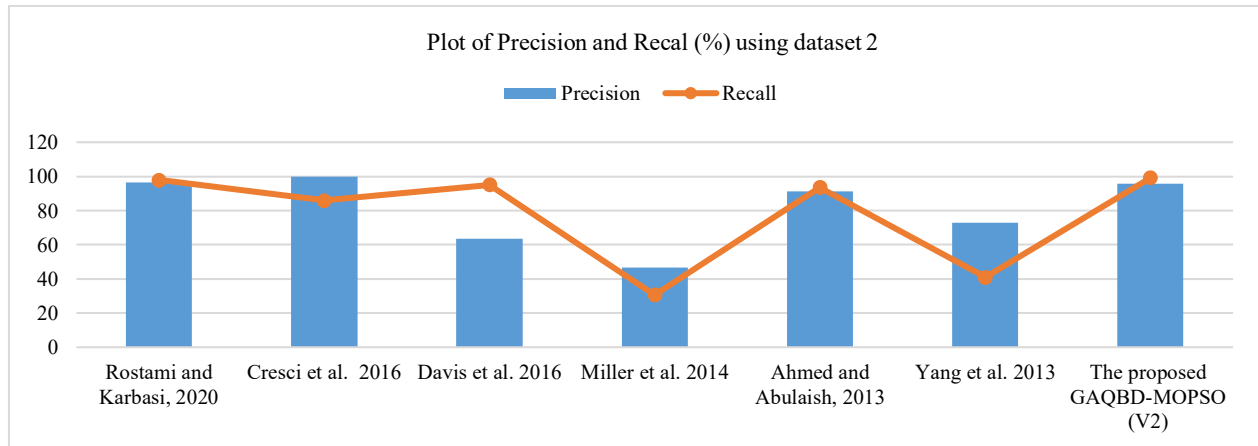
feature vector for a set of unlabelled samples in the dataset. Figures 15 and 16 show that the clustering algorithm proposed in [24] doesn't achieve good recall in datasets 1 and 2. So it's very difficult to detect spambots from data streams. Ahmed and Abulaish [19], proposed a graph clustering algorithm based on Markov Clustering Algorithm (MCL) [55]. However, Cresci *et al.* [10] have replaced the MCL algorithm with the Fastgreedy community detection algorithm [56], thus avoiding the problem of identifying two distinct clusters.

Inspired by the biological DNA sequence, Cresci *et al.* [22] proposed a bio-inspired model called digital DNA system aimed at recognising the behaviour of online users. The digital DNA was expressed by a string encoding each user's behaviour. Then, the Longest Common Substring (LCS) measurement was used to determine the anomalous similarities between the sequences, and the longest DNA sequences were labelled as spambot accounts. From Figures 15 and 16, the recall rate for datasets 1 and 2 is unbalanced, leading to a reduction in the performance criteria. Furthermore, Rostami and Karbasi [17] proposed a multi-objective feature selection approach to select a stable subset of features based on the highest relation to the target class and the least redundancy among the features using the Minimum Redundancy – Maximum Relevance algorithm (mRMR) [18]. In this study, the QB-MOPSO algorithm for detecting fake accounts on Twitter is presented and denoted the Neuro- QB-MOPSO system. Figures 15 and 16 illustrate the superiority of the algorithm GAQB-MOPSO -V2 compared to other methods based on the precision and recall criteria. Neuro-QB-MOPSO is a supervised method that is inspired by the standard quantum PSO algorithm. It can detect human and fake accounts on Twitter much more accurately than other methods.



**Figure 14.** The Precision and Recall Rate of GAQB-MOPSO (V2) Compared with state-of-the-art Methods using Dataset 1.





**Figure 15.** The Precision and Recall Rate of GAQB-MOPSO (V2) Compared with state-of-the-art Methods using Dataset 2.

## 5. Conclusion

This paper proposed a new Quantum Beta-behaved Multi-Objective Particle Swarm Optimization (QB-MOPSO) algorithm, comprising six quantum variants with different beta-profiles. The QB-MOPSO system was used for pertinent feature selection to detect fake accounts on Twitter. The main goal was to minimize both the features' dimensionality and the classification error rate. The six variants of the QB-MOPSO approach were proposed with two optimization profiles, the first was for exploration using a quantum-behaved MOPSO, and the second was for exploitation phase using a Beta-behaved MOPSO. Both profiles were assumed over new mathematical rules to optimize and update the velocities and the positions of particles in the search space. At each iteration, binary encoding is fixed using the sigmoid function. Therefore, the bit '1' indicates a selected feature and '0' otherwise. The proposed system was tested on the two benchmark Twitter datasets and achieved excellent results compared with state-of-the-art methods. The GAQB-MOPSO (V2) system was found to achieve an accuracy rate of 99.19% on dataset 1 and 97.52% for dataset 2. For future work, we will address the challenge of online feature selection to predict online fake accounts on OSNs, considering the stability of the feature subset. Also, a new investigation will be proposed for fake news detection on Twitter.

## Acknowledgment

The research leading to these results has received funding from the Ministry of Higher Education and Scientific Research of Tunisia under the grant agreement number LR11ES48.

## Statements and Declarations

### Competing Interests

The authors declare that we do not have any commercial or associative interest that represents a conflict of interest in connection with the work submitted

### Authors contribution statement

**Ahlem Aboud:** Conceptualization, Methodology, Software, Validation, Writing – Original Draft, Writing Review & Editing, Visualization, Formal analysis, Experimental investigations. **Nizar Rokbani:** Supervision, Conceptualization, Methodology, Validation, Formal analysis, Writing – Original Draft. **Seyedali Mirjalili:** Formal analysis, Conceptualization, Methodology, Validation. **Amir Hussain:** Supervision, Paper arrangements, Formal analysis, Validation. **Habib Chabchoub:** Supervision, Paper arrangements, Formal analysis, Validation. **Adel M. Alimi:** Supervision, Project administration, Conceptualization, Validation, Writing – Original Draft.

### Ethical and informed consent for data used

Not applicable

### Data availability and access

The authors declare that all data supporting the findings of this study are available within the article.

### References

- [1] S. B.-T. I. Encyclopedia and undefined 2003, “Personalization and customization technologies,” *books.google.com*, Accessed: Feb. 01, 2022. [Online]. Available: <https://books.google.com/books?hl=fr&lr=&id=wsh3f0hyI8C&oi=fnd&pg=PA51&dq=Personalization+and+Customization+Technologies&ots=-j-m2hYRy2&sig=fQuECXOHzApigtn9CWKdiRMro>
- [2] S. Kanoje, S. Girase, and D. Mukhopadhyay, “User Profiling Trends, Techniques and Applications,” *undefined*, vol. 1, no. 1, 2015.
- [3] Z. Alom, B. Carminati, and E. Ferrari, “Detecting spam accounts on Twitter,” *Proceedings of the 2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2018*, pp. 1191–1198, Oct. 2018, doi: 10.1109/ASONAM.2018.8508495.
- [4] A. Azab, A. M. Idrees, M. A. Mahmoud, and H. Hefny, “Fake Account Detection in Twitter Based on Minimum Weighted Feature set,” *undefined*, 2015.
- [5] J. Castellini, V. Poggioni, and G. Sorbi, “Fake twitter followers detection by denoising autoencoder,” *Proceedings - 2017 IEEE/WIC/ACM International Conference on Web Intelligence, WI 2017*, pp. 195–202, Aug. 2017, doi: 10.1145/3106426.3106489.
- [6] S. R. Sahoo and B. B. Gupta, “Multiple features based approach for automatic fake news detection on social networks using deep learning,” *Appl Soft Comput*, vol. 100, p. 106983, Mar. 2021, doi: 10.1016/J.ASOC.2020.106983.
- [7] F. C. Akyon and M. Esat Kalfaoglu, “Instagram Fake and Automated Account Detection,” *Proceedings - 2019 Innovations in Intelligent Systems and Applications Conference, ASYU 2019*, Oct. 2019, doi: 10.1109/ASYU48272.2019.8946437.
- [8] J. Miao and L. Niu, “A Survey on Feature Selection,” *Procedia Comput Sci*, vol. 91, pp. 919–926, Jan. 2016, doi: 10.1016/J.PROCS.2016.07.111.
- [9] B. Xue, M. Zhang, and W. N. Browne, “Particle swarm optimization for feature selection in classification: A multi-objective approach,” *IEEE Trans Cybern*, vol. 43, no. 6, pp. 1656–1671, Dec. 2013, doi: 10.1109/TSMCB.2012.2227469.
- [10] S. Cresci, R. di Pietro, M. Petrocchi, A. Spognardi, and M. Tesconi, “The paradigm-shift of social spambots: Evidence, theories, and tools for the arms race,” *26th International World Wide Web Conference 2017, WWW 2017 Companion*, pp. 963–972, Jan. 2017, doi: 10.1145/3041021.3055135.
- [11] J. Sun, W. Xu, and B. Feng, “A global search strategy of Quantum-behaved Particle Swarm Optimization,” *2004 IEEE Conference on Cybernetics and Intelligent Systems*, pp. 111–116, 2004, doi: 10.1109/ICCIS.2004.1460396.
- [12] J. Sun, W. Fang, V. Palade, X. Wu, and W. Xu, “Quantum-behaved particle swarm optimization with Gaussian distributed local attractor point,” *Appl Math Comput*, vol. 218, no. 7, pp. 3763–3775, Dec. 2011, doi: 10.1016/J.AMC.2011.09.021.
- [13] A. M. Alimi, “BETA NEURO-FUZZY SYSTEMS,” vol. 1, no. 1, pp. 23–41, 2003.

- [14] F. Morstatter, L. Wu, T. H. Nazer, K. M. Carley, and H. Liu, "A new approach to bot detection: Striking the balance between precision and recall," *Proceedings of the 2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2016*, pp. 533–540, Nov. 2016, doi: 10.1109/ASONAM.2016.7752287.
- [15] K. Sutha and J. Tamilselvi, "A Review of Feature Selection Algorithms for Data Mining Techniques," *undefined*, 2015.
- [16] J. Tang, S. Alelyani, H. L.-D. classification: A. and, and undefined 2014, "Feature selection for classification: A review," *cvs.edu.in*, Accessed: Feb. 02, 2022. [Online]. Available: [http://www.cvs.edu.in/upload/feature\\_selection\\_for\\_classification.pdf](http://www.cvs.edu.in/upload/feature_selection_for_classification.pdf)
- [17] R. R. Rostami and S. Karbasi, "Detecting fake accounts on twitter social network using multi-objective hybrid feature selection approach," *Webology*, vol. 17, no. 1, Jun. 2020, doi: 10.14704/WEB/V17I1/A204.
- [18] C. Ding and H. Peng, "Minimum redundancy feature selection from microarray gene expression data," *Proceedings of the 2003 IEEE Bioinformatics Conference, CSB 2003*, pp. 523–528, 2003, doi: 10.1109/CSB.2003.1227396.
- [19] F. Ahmed and M. Abulaish, "A generic statistical approach for spam detection in Online Social Networks," *Comput Commun*, vol. 36, no. 10–11, pp. 1120–1129, Jun. 2013, doi: 10.1016/J.COMCOM.2013.04.004.
- [20] H. Peng, F. Long, and C. Ding, "Feature selection based on mutual information: Criteria of Max-Dependency, Max-Relevance, and Min-Redundancy," *IEEE Trans Pattern Anal Mach Intell*, vol. 27, no. 8, pp. 1226–1238, Aug. 2005, doi: 10.1109/TPAMI.2005.159.
- [21] C. A. Davis, O. Varol, E. Ferrara, A. Flammini, and F. Menczer, "BotOrNot," pp. 273–274, 2016, doi: 10.1145/2872518.2889302.
- [22] S. Cresci, R. di Pietro, M. Petrocchi, A. Spognardi, and M. Tesconi, "DNA-Inspired Online Behavioral Modeling and Its Application to Spambot Detection," *IEEE Intell Syst*, vol. 31, no. 5, pp. 58–64, Sep. 2016, doi: 10.1109/MIS.2016.29.
- [23] C. Yang, R. Harkreader, and G. Gu, "Empirical evaluation and new design for fighting evolving twitter spammers," *IEEE Transactions on Information Forensics and Security*, vol. 8, no. 8, pp. 1280–1293, 2013, doi: 10.1109/TIFS.2013.2267732.
- [24] Z. Miller, B. Dickinson, W. Deitrick, W. Hu, and A. H. Wang, "Twitter spammer detection using data stream clustering," *Inf Sci (N Y)*, vol. 260, pp. 64–73, Mar. 2014, doi: 10.1016/J.INS.2013.11.016.
- [25] S. Lee, B. Schowe, and V. Sivakumar, "Feature Selection for High-Dimensional Data with RapidMiner," *undefined*, 2012, doi: 10.17877/DE290R-14289.
- [26] A. Aboud, R. Fdhila, and A. M. Alimi, "Dynamic Multi Objective Particle Swarm Optimization Based on a New Environment Change Detection Strategy," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, Nov. 2017, vol. 10637 LNCS, pp. 258–268. doi: 10.1007/978-3-319-70093-9\_27.
- [27] A. Aboud, R. Fdhila, and A. M. Alimi, "MOPSO for dynamic feature selection problem based big data fusion," in *2016 IEEE International Conference on Systems, Man, and Cybernetics, SMC 2016 - Conference Proceedings*, Feb. 2017, pp. 3918–3923. doi: 10.1109/SMC.2016.7844846.
- [28] A. Aboud, R. Fdhila, A. Hussain, and A. M. Alimi, "A novel Dynamic Pareto bi-level Multi-Objective Particle Swarm Optimization (DPb-MOPSO) algorithm," *TechRxiv*, Dec. 2020. doi: 10.36227/TECHRXIV.13325354.V1.
- [29] A. Aboud, N. Rokbani, S. Mirjalili, and A. Alimi, "A Distributed Bi-behaviors Crow Search Algorithm for Dynamic Multi-Objective Optimization and Many-Objective Optimization," Sep. 2021, doi: 10.36227/TECHRXIV.16607858.V2.
- [30] A. Aboud *et al.*, "A Distributed Multifactorial Particle Swarm Optimization Approach," Dec. 2021, doi: 10.36227/TECHRXIV.17260040.V1.
- [31] A. Aboud *et al.*, "DPb-MOPSO: A Novel Dynamic Pareto bi-level Multi-Objective Particle Swarm Optimization Algorithm," Dec. 2021, doi: 10.36227/TECHRXIV.17207576.V1.
- [32] A. E. Ezugwu, J. O. Agushaka, L. Abualigah, S. Mirjalili, and A. H. Gandomi, "Prairie Dog Optimization Algorithm," *Neural Comput Appl*, pp. 1–49, Jul. 2022, doi: 10.1007/S00521-022-07530-9/FIGURES/31.
- [33] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, "A fast and elitist multiobjective genetic algorithm: NSGA-II," *IEEE Transactions on Evolutionary Computation*, vol. 6, no. 2, pp. 182–197, Apr. 2002, doi: 10.1109/4235.996017.
- [34] J. Kennedy and R. Eberhart, "Particle swarm optimization," in *Proceedings of ICNN'95 - International Conference on Neural Networks*, vol. 4, pp. 1942–1948. doi: 10.1109/ICNN.1995.488968.

- [35] D. A. Savic, G. A. Walters, and J. W. Davidson, "A Genetic Programming Approach to Rainfall-Runoff Modelling," *Water Resources Management* 1999 13:3, vol. 13, no. 3, pp. 219–231, 1999, doi: 10.1023/A:1008132509589.
- [36] M. Dorigo, M. Birattari, and T. Stutzle, "Ant colony optimization," *IEEE Comput Intell Mag*, vol. 1, no. 4, pp. 28–39, Nov. 2006, doi: 10.1109/MCI.2006.329691.
- [37] B. Xue, M. Zhang, ... W. B.-I. T. on, and undefined 2015, "A survey on evolutionary computation approaches to feature selection," *ieeexplore.ieee.org*, Accessed: Feb. 02, 2022. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/7339682/>
- [38] A. W. Whitney, "A Direct Method of Nonparametric Measurement Selection," *IEEE Transactions on Computers*, vol. C-20, no. 9, pp. 1100–1103, 1971, doi: 10.1109/T-C.1971.223410.
- [39] T. Marill and D. M. Green, "On the Effectiveness of Receptors in Recognition Systems," *IEEE Trans Inf Theory*, vol. 9, no. 1, pp. 11–17, 1963, doi: 10.1109/TIT.1963.1057810.
- [40] Y.-S. Jeong, K. S. Shin, and M. K. Jeong, "An evolutionary algorithm with the partial sequential forward floating search mutation for large-scale feature selection problems," *Journal of the Operational Research Society* 2014 66:4, vol. 66, no. 4, pp. 529–538, Mar. 2014, doi: 10.1057/JORS.2013.72.
- [41] S. M. Winkler, M. Affenzeller, W. Jacak, and H. Stekel, "Identification of cancer diagnosis estimation models using evolutionary algorithms: A case study for breast cancer, melanoma, and cancer in the respiratory system," *Genetic and Evolutionary Computation Conference, GECCO'11 - Companion Publication*, pp. 503–510, 2011, doi: 10.1145/2001858.2002040.
- [42] Y. Zhang, D. Gong, Y. Hu, and W. Zhang, "Feature selection algorithm based on bare bones particle swarm optimization," *Neurocomputing*, vol. 148, pp. 150–157, Jan. 2015, doi: 10.1016/J.NEUCOM.2012.09.049.
- [43] B. Xue, M. Zhang, and W. N. Browne, "Particle swarm optimisation for feature selection in classification: Novel initialisation and updating mechanisms," *Appl Soft Comput*, vol. 18, pp. 261–276, May 2014, doi: 10.1016/J.ASOC.2013.09.018.
- [44] S. M. Vieira, J. M. C. Sousa, and T. A. Runkler, "Two cooperative ant colonies for feature selection using fuzzy models," *Expert Syst Appl*, vol. 37, no. 4, pp. 2714–2723, Apr. 2010, doi: 10.1016/J.ESWA.2009.08.026.
- [45] K. K. Bharti and S. Pandey, "Fake account detection in twitter using logistic regression with particle swarm optimization," *Soft comput*, vol. 25, no. 16, pp. 11333–11345, Aug. 2021, doi: 10.1007/S00500-021-05930-Y/FIGURES/5.
- [46] G. Lingam, R. R. Rout, and D. V. L. N. Somayajulu, "Deep Q-Learning and Particle Swarm Optimization for Bot Detection in Online Social Networks," *2019 10th International Conference on Computing, Communication and Networking Technologies, ICCCNT 2019*, Jul. 2019, doi: 10.1109/ICCCNT45670.2019.8944493.
- [47] J. Sun, B. Feng, and W. Xu, "Particle swarm optimization with particles having quantum behavior," *Proceedings of the 2004 Congress on Evolutionary Computation, CEC2004*, vol. 1, pp. 325–331, 2004, doi: 10.1109/CEC.2004.1330875.
- [48] C. A. Coello Coello and M. S. Lechuga, "MOPSO: A proposal for multiple objective particle swarm optimization," *Proceedings of the 2002 Congress on Evolutionary Computation, CEC 2002*, vol. 2, pp. 1051–1056, 2002, doi: 10.1109/CEC.2002.1004388.
- [49] Y. Qi, Q. Zhang, X. Ma, Y. Quan, and Q. Miao, "Utopian point based decomposition for multi-objective optimization problems with complicated Pareto fronts," *Appl Soft Comput*, vol. 61, pp. 844–859, Dec. 2017, doi: 10.1016/J.ASOC.2017.08.036.
- [50] A. Farkiya, P. Saini, S. Sinha, and S. Desai, "Natural Language Processing using NLTK and WordNet," 2015, Accessed: Feb. 02, 2022. [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.735.2951&rep=rep1&type=pdf>
- [51] H. Kaur, V. Mangat, and Nidhi, "A survey of sentiment analysis techniques," *Proceedings of the International Conference on IoT in Social, Mobile, Analytics and Cloud, I-SMAC 2017*, pp. 921–925, Oct. 2017, doi: 10.1109/I-SMAC.2017.8058315.
- [52] Genichi. Taguchi, Rajesh. Jugulum, and Shin. Taguchi, "Computer-Based Robust Engineering : an Essential for DFSS.," p. 240, 2004.
- [53] L. Jinxian and L. Hui, "A density-based clustering over evolving heterogeneous data stream," *2009 Second ISECS International Colloquium on Computing, Communication, Control, and Management, CCCM 2009*, vol. 4, pp. 275–277, 2009, doi: 10.1109/CCCM.2009.5267735.
- [54] M. R. Ackermann, M. Märtens, C. Raupach, K. Swierkot, C. Lammersen, and C. Sohler, "StreamKM++," *Journal of Experimental Algorithmics (JEA)*, vol. 17, May 2012, doi: 10.1145/2133803.2184450.

- [55] S. van Dongen, “Graph Clustering Via a Discrete Uncoupling Process,” <http://dx.doi.org/10.1137/040608635>, vol. 30, no. 1, pp. 121–141, Feb. 2008, doi: 10.1137/040608635.
- [56] A. Clauset, M. E. J. Newman, and C. Moore, “Finding community structure in very large networks,” *Phys Rev E Stat Phys Plasmas Fluids Relat Interdiscip Topics*, vol. 70, no. 6, p. 6, Dec. 2004, doi: 10.1103/PHYSREVE.70.066111/FIGURES/3/MEDIUM.