

Multi-objective optimization of base classifiers in StackingC by NSGA-II for intrusion detection

Michael Milliken*, Yaxin Bi[†], Leo Galway[†] and Glenn Hawe[†]

School of Computing and Mathematics, Ulster University, Belfast, United Kingdom

*Email: milliken-m@email.ulster.ac.uk

[†]Email: {y.bi, l.galway, gi.hawe}@ulster.ac.uk

Abstract—Multiple Classifier Systems are often found to improve results of intrusion detection by combining a set of classifier decisions where single classifiers may not achieve the same level of detection. However not every set of classifiers is more able, therefore selection of more capable sets is required. A misclassification is a false positive or negative instance; a set of classifiers may produce one more than the other. An optimal set of classifiers is required to reduce both, thus treating them as individual objectives allows a balance to be found. The aim of this work is the selection of optimal sets of base level classifiers using an evolutionary computation approach. A comparative analysis is made of the performance of the generated ensembles against the individual base level classifiers, it is shown that optimal ensembles can be found to perform better than a majority of individuals.

I. INTRODUCTION

An ongoing challenge, for security experts and businesses, in the defence and stability of modern networks is the detection of intrusions that either illicitly retrieve data/information, or prevent legitimate access [1]. Given the wide varieties of existing intrusions, one key difficulty is providing accurate detection. Intrusion detection depends largely on analysis of attack vectors and purposes. Attack methods and perceived purposes may be useful to determine if they are different enough from normal traffic to accurately state them to be attacks.

Systems put in place to detect intrusions are known as Network Intrusion Detection Systems (NIDS). An NIDS may function with respect to historical or real time network traffic [2]. In either case, they are utilized to provide some method of detection through analysing all or a specified set of features [2], [3]. A popular approach to NIDS over the last decade has been the use of multi-layered or tiered approaches [4], [5], [6], which incorporate multiple methods within a hierarchy. Accuracy of intrusions depends upon the algorithm employed, when a single algorithm is not adequate, multiple algorithms may be required, determination of how many and which algorithms to employ requires further consideration.

The focus of this paper concerns the use of a Multi-objective Genetic Algorithm (GA) to determine Pareto-optimal ensembles of Base Level Classifiers (BLC) for the detection of intrusions. The hypothesis is that the evolution of an ensemble of BLCs may result in better performance, to some degree, than individual BLCs. To explore this, we measured performances of a set of supervised Machine Learning (ML) algorithms, previously used with a GA [7], on two datasets

used for intrusion detection research, notably the NSL-KDD [8] and ISCX2012 [9] datasets. The developed work is the use of an evolutionary approach to sets of ML algorithms in the application domain of network intrusion detection. The work is differentiated by the use of the specific ML algorithms for an efficient ensemble StackingC for the particular use of network intrusion detection with more recent and large dataset subsets. In particular, investigating performance and determining applicability of produced ensembles per dataset as well as existence of general ensemble across each.

The remainder of this paper is outlined as follows: Section II describes some background on GA, ML algorithms and existing NIDS. Following this, Section III, describes the methodology employed for the experiments discussed herein, followed by a presentation of the results obtained in Section IV. Finally, in Section V a brief discussion of the results is given along with a subsequent conclusion in Section VI.

II. BACKGROUND

As ML research and development continues more potential algorithms are produced. For any given situation one algorithm may perform as well as another algorithm, in this case the need for both algorithms may become less likely, at least for a certain series of instances [7]. Each algorithm may require a certain amount of time for training and classification, the overall amount of time may not simply be linear and predictable. Therefore, implementation of all algorithms may not be desirable or adequate for the problem of accurate classification.

In general, the set of optimal algorithms for a given problem is known to be a subset of all possible algorithms [10]. Thus, selection of an optimal set requires a further optimization process, such as a GA which follows the theory of evolution, producing optimal sets of classifiers with objective evaluation criteria.

NIDS can be evaluated using various different measures. For example, F-Measure (FM), Detection Rate, totals of False Positive (FP) or False Negative (FN), or rates of FP or FN classification instances may be employed. Within this paper, FP is taken to be the erroneous classification of an instance as an intrusion. Conversely, a FN is taken to be the erroneous classification of an instance as normal. With NIDS a balance must be found between FPs and FNs, as these two objectives generally compete. One way to achieve this balance

is to identify the Pareto-optimal trade-off between FP and FN using Multi-objective optimization algorithms. A noted limiting factor of NIDS is the reduction of FPs [11], [12].

A. Ensemble Algorithms

Outputs from multiple BLCs can be combined into a single output, this approach may be referred to as an *Ensemble* or *Multiple Classifier System (MCS)*. Popular conventional Ensemble methods include Bagging, Boosting and Stacking. Bagging uses instances of a BLC and a replicated dataset with differences introducing variations, improving average classification result [13]. By contrast, Boosting uses BLCs sequentially over one dataset. Improvements are induced at each subsequent BLC, weighting changes of misclassified instances, improving overall classification [14]. Similar to Bagging's use of multiple BLCs and Boosting's use of one dataset, Stacking uses BLCs in parallel over one dataset [15]. Each BLC predicts class probabilities as input for a regression model per class. Each regression model classifies each instance [16]. Using class probabilities rather than predictions improve overall performance. A Stacking variation, StackingC [17], improved efficiency by using only the probability of a specific class for each linear model for that class, rather than considering the other class probabilities as well [17].

B. Ensemble-based NIDS

Within research literature, some NIDS have utilized Ensemble approaches, improving upon performance of a single BLC.

Octopus-IIDS [5], described as an Intelligent IDS, implements a Kohonen network to split data into attack classes, then SVMs to reclassify instances as *Normal* or *Attack*.

Hidden Markov Models with Payl (HMMPayl) [18] performs analysis on payloads, as well as using Hidden Markov Models (HMMs) initially, performing multiple initial classifications, then forming a final classification from the initial classification results.

Work described in [19] forms clusters with K-Means Clustering (KMC), after which Naïve Bayes (NB) corrects data previously misclassified by the clustering. This approach increases classification results as well as efficiency by grouping the data.

C. Genetic Algorithms

A GA is an evolutionary algorithm that mimics the process of natural selection to evolve a population of candidate solutions to an optimization problem [20], [21]. The evolutionary aspect of the algorithm is related to modelling evolutionary processes whereby two chromosomes are 'mated', producing 'child' chromosomes based on 'parent' chromosome representations. Improvements can arise by selection and crossover (i.e. 'mating') of the fittest chromosomes. Mutations may also be applied, whereby a gene of a chromosome is randomly modified to introduce diversity within the population of chromosomes. Each chromosome represents a potential solution to the problem being solved.

Objective functions of a GA generate values utilized to compare performance of decoded chromosomes, such that

the optimal or dominating chromosome may be identified. Values may be maximized or minimized depending upon need; for instance, maximising accuracy and minimizing error. GAs with two or more objectives are Multi-objective GAs (MOGA). In Multi-objective optimization, often the multiple objectives compete. For competing objectives there exists no 'global' solution, therefore a MOGA must use the concept of dominance to compare and rank solutions. A solution x_1 is said to dominate another solution x_2 if and only if x_1 is strictly better than x_2 in one objective, and no worse in all the others. Solutions that are not dominated by any other solution represent the best possible trade-off in objectives, and are called Pareto-optimal. The goal of a MOGA is to identify a diverse set of Pareto-optimal solutions.

One of the most popular MOGA is Non-dominated Sorting Genetic Algorithm II (NSGA-II) [22], an improved version of NSGA. NSGA-II creates an initial population of chromosomes and evaluates those chromosomes using a defined problem. Chromosomes in pairs create a child chromosome population using crossover and mutation, which is then evaluated. Evaluated chromosomes are all compared to each other, providing the domination information of each chromosome, i.e. which chromosomes dominated and which were dominated by an individual chromosome. Non-dominated chromosomes form the first or "0" subset; these Pareto-optimal chromosomes are removed from any dominating set of remaining chromosomes. This process continues, finding all possible and existing Pareto fronts for the problem until all chromosomes have been assigned to a front.

NSGA-III [23], is an improved version of NSGA-II, able to more efficiently evolve solutions where problems include a large number of objectives.

MOEA/D [24], is a Multi-objective Evolutionary Algorithm that decomposes optimization problems into sub-problems, at each generation the best solution of each sub-problem forms the new generation, in this way the complexity of the overall problem is reduced.

D. Intrusion-based NIDS using GAs

GAs are also used to choose an optimal set of variables. In [25] a GA has been implemented and used for the detection of intrusions from the KDD'99 dataset. A multi-class approach was taken, performing well on Denial of Service (DoS) data, however not as well on Normal data, although the performance on DoS and User to Root data is stated as better than that of the KDD'99 competition winner.

In [26] a set of features and the parameters of SVMs kernel function is evolved to optimize the detection rate of the resultant SVM. The proposed system is said to sometimes outperform the KDD'99 competition winner. Similarly in [27], a GA named Archive-based Micro Genetic Algorithm 2 (AMGA2) is used to optimize features, in this case those of NB. The results from the evolved NB are compared against other methods using NB as the BLC. Experiments are carried out on both KDD'99 and the more recent ISCX2012

datasets, although low numbers of instances are selected from ISCX2012.

III. METHODOLOGY

The focus of the work presented in this paper is investigation and determination of optimal sets of BLCs, using NSGA-II, for use in an ensemble, StackingC, employed for the detection of network intrusions.

BLCs and an efficient ensemble method are implemented, as described in [7]. A particular distinction is the difference of domain upon which the set of BLCs are used, in this case a NIDS. Additionally, in comparison with the current literature, a more recent dataset, ISCX2012, is also employed with large number of instances, facilitating investigation into the effects of large scales of data with StackingC in an evolutionary context.

To simplify selection and implementation of algorithms for the ensemble, algorithms within Weka are used to form BLCs and the Meta Level Classifier (MLC), e.g. Multi-response Linear Regression. In addition, JMetal [28], an "object-oriented Java-based framework for multi-objective optimization with metaheuristics", has been utilized for encoding and decoding chromosomes for the GA, providing an initial population of chromosomes and assisting in catching potential issues with child chromosomes; ensuring a minimum number of BLCs in initial populations and handling when chromosomes decode to an empty set.

Although the most recent variation of NSGA, NSGA-III, is available, NSGA-II provides a robust and popular approach that is adequate for optimization of two potentially competing objectives. Development of a bespoke GA was beyond the scope of this paper, thus a generic implementation of NSGA-II was implemented.

A. StackingC with pre-trained BLCs

StackingC required a number of changes to provide the efficiency found when using pre-trained BLCs. Objects and methods were added to Stacking and StackingC (the latter depending upon the former), maintaining trained BLCs, indices selecting required BLCs for the current chromosome and the meta data with which to train each MLC.

With the existence of a possibly unlimited number of BLCs for use in an ensemble, StackingC needs to be more efficient. When StackingC is used over multiple iterations of a set of BLCs, each set would require training on the full dataset for each iteration. Considering all BLCs and the amount of training data, the training time required for each ensemble is potentially large. Thus due to a combination of, firstly, size of current and future BLC sets, secondly, training datasets and, thirdly, time required to train individual BLCs, a change to only training once was required.

B. NSGA-II using BLC evaluations

Binary encoding of the set of BLCs was utilised, whereby the presence or absence of a BLC is represented by a 1 or 0 respectively. An implementation of NSGA-II for binary

encodings was utilized. Subsequently, creation and evolution of chromosomes provide sets of BLCs, where these sets are subsets of the entire possible set of BLCs. Comparisons of chromosomes involved minimization of both FP and FN values. It is assumed that, two sets of BLCs that perform well may be evolved and produce at least one set of BLCs that can perform better than and is different to the prior two sets.

C. Evolving sets of BLCs for StackingC

The use of a GA would facilitate the process of evolving a chromosome evaluated on pre-selected objectives, evaluation provided by the more efficient StackingC. Within this section such a system is described, illustrated in Fig.1. The illustration provides 10 stages in total, the methodology of the experiment from an initial StackingC and population generation followed by descriptions of the evolution of the population.

An initial iteration of StackingC is run on training data, training the MLC but also producing a set of Trained BLCs from the full set of BLCs and the Meta Data, as shown in stages 1-3 of Fig.1 respectively. Trained BLCs and Meta Data are represented in stages 2 and 3 respectively. Subsequently, as illustrated in stage 4 of Fig.1, NSGA-II through the use of a problem description creates multiple solutions forming the initial population of chromosomes.

Following the formation of an initial population illustrated in stage 4, in stage 5 of Fig.1 the GA decodes a chromosome from the population to produce an array of indices indicating which BLC should be present in a set. This array is used against the full set of Trained BLCs to provide the set of Reduced BLCs, given by stage 6 of Fig.1.

Both the array and Reduced BLCs are passed to a new StackingC instance, given by stage 7 of Fig.1, to set the required BLCs and the indices of data from the Meta Data initially used to train the MLC, as shown in stage 8 of Fig.1.

The combination of these three objects, the Reduced BLCs, the array and the Meta Data, provides the creation of a StackingC model without retraining the whole set of BLCs or recreating Meta Data with which to train the MLCs, as shown in stage 8 of Fig.1. Within this new StackingC instance, Meta Data is reduced using the array to form Reduced Meta Data. This reduction is required as, similarly with Reduced BLCs, the new StackingC instance will only require specific parts of Meta Data to produce Trained MLCs. These specific parts being representative of the ensemble provided by NSGA-II.

Testing Data passed to the Reduced BLCs is used to produce Testing Meta Data for the testing of Trained MLCs. Subsequent Classification and evaluation of the ensemble, as described by the chromosome, are used to update the solution's objective values (Update Solution) and subsequently the solution in the population (Update Population), as shown in stage 9 of Fig.1. The process illustrated in stages 5-9 of Fig.1 is repeated for each chromosome in the initial population, after which the process illustrated in stage 10 is performed.

With the population evaluated, if the stopping condition is not met, which in this case is the evaluation of 30 generations of the population, a child population is evolved. Each child

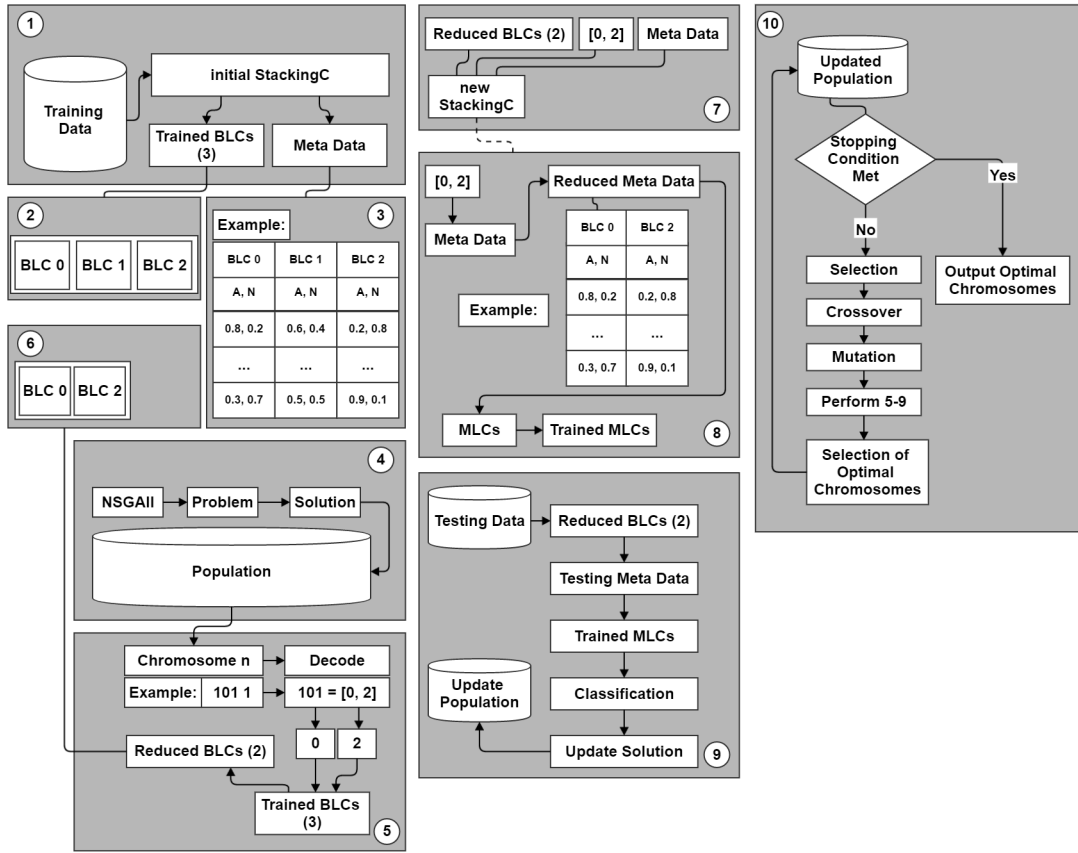


Fig. 1. Illustration of efficient StackingC training and subsequent evolution by NSGA-II

population chromosome is subject to processes of stages 5-9 of Fig.1. Optimal chromosomes from parent and child set are found by measure of domination, as described previously, to form a new Updated Population of parent chromosomes, as shown in stage 10 of Fig.1. The process repeats until the specified number of generations has been passed.

D. Datasets

Work presented in this paper employs two primary datasets: NSL-KDD and ISCX2012. NSL-KDD has been pre-processed into training and testing datasets, the proportion of which is presented in Table I. By contrast, ISCX2012 has not been pre-processed, hence required training and testing datasets to be subsequently generated. Consequently, the ISCX2012 dataset was split in two ways: UUJN - using only the Distributed Denial of Service (DDoS) scenario from the dataset, with a similar distribution of classes as NSL-KDD; and Yassin2013 - using an approximation of the dataset used by [19] with a higher number of attack instances in the training set. Table I gives brief description of all the dataset variation datasets utilized within the experiments.

IV. RESULTS

A number of experiments were performed to determine the optimal sets of BLCs for StackingC, with focus on detecting

network intrusions. Each dataset was used with the full set of BLCs during the experiments. Adjustments to the full set of BLCs were required to allow results from each dataset to be obtained. This was due to different feature types, number of instances as well as BLC performance per dataset.

One issue faced was time complexity of training; When trained on NSL-KDD, the full set of BLCs were able to produce results within an adequate time frame. However, UUJN and Yassin2013 were found to require extended periods of time for training. Some BLCs required scores of hours before finishing and in some cases did not finish. Therefore, results from only a subset of the BLCs for these two datasets, as indicated in Table II, were viable. Consequently, BLCs that experienced extended training times were removed from the full set and a new subset/full set was produced. This new set experienced no such issues with UUJN, however still some issues were found with Yassin2013 for one BLC. In response, again, the BLC causing issues was removed and a new subset/full set was produced.

Thus three BLC set variations were of appropriate use for each dataset with some overlap between which datasets and which sets of BLCs provided results. Positions of each BLC for each BLC set variation in which they were viable are presented in Table II under the column **Base Level Classifier**

TABLE I
CLASS DISTRIBUTIONS OF DATASETS USED FOR EXPERIMENTS

Dataset	Features	Training			Testing		
		Normal	Attack	Total	Normal	Attack	Total
NSL-KDD	42	67,343	58,630	125,973	9,711	12,833	22,544
UUJN	16	35,223	30,666	65,889	4,571	6,040	10,611
Yassin2013	15	63,765	8,968	72,733	19,115	37,159	56,274

TABLE II
PRESENCE OF BASE LEVEL CLASSIFIERS FROM WEKA

Base Level Classifier		Base Level Classifier Set Variation		
Name	Type	One	Two	Three
<i>NaïveBayesUpdateable</i>	bayes	10000000000000	10000000000	1000000000
<i>PART</i>	rules	01000000000000	01000000000	0100000000
<i>J48 (pruned)</i>	trees	00100000000000	00100000000	0010000000
<i>J48 (unpruned)</i>	trees	00010000000000	00010000000	0001000000
<i>DecisionStump</i>	trees	00001000000000	00001000000	0000100000
<i>DecisionTable</i>	rules	00000100000000	00000100000	0000010000
<i>ClassificationViaRegression</i>	meta	00000010000000	00000010000	0000001000
<i>RandomForest</i>	trees	00000001000000		
<i>RandomTree</i>	trees	00000000010000		
<i>VFI</i>	misc	00000000010000	00000010000	0000001000
<i>ConjunctiveRule</i>	rules	00000000001000	00000001000	0000000100
<i>JRip</i>	rules	00000000000010	00000000100	0000000010
<i>NNge</i>	rules	000000000000010	00000000010	0000000001
<i>HyperPipes</i>	misc	000000000000001	00000000001	0000000001
Total Classifiers		14	11	10

Set Variation. The three BLC set variations are thus named as *One*, *Two* and *Three* or referred to respectively as first, second and third BLC set variation. NSL-KDD produced results with each BLC set variation, UUN produced results with the second and third BLC set variations and Yassin2013 only produced results with the third BLC set variation. This is shown in Tables II, III and IV.

A. On NSL-KDD Results

The results of each BLC set variation, as given in Table III, show that overall there is a single data point that exists as an optimal set of FP and FN values, achievable by multiple unique ensembles. This is true across each BLC set variation where the FP and FN values obtained remain the same. All ensembles across NSL-KDD achieve 248 FPs and 3263 FNs, while some individual BLCs achieve similar results. From Table IV it is seen that ensembles perform best regarding FN where in all but one case the individuals achieve around 700 more FNs.

Considering the results obtained using NSL-KDD, shown in Table III and Table IV, it can be determined that ensembles outperform the majority of BLCs, in terms of FP as well as FN. An exception, however, is found with the HyperPipes BLC, which outperforms any of the ensembles in terms of FP, achieving 56 FP compared to 248 FP achieved by any ensemble.

The data points from BLCs and optimal ensembles are presented in Fig.2. It is evident from positions of BLCs in comparison to the optimal ensemble that the ensemble outperforms the BLCs.

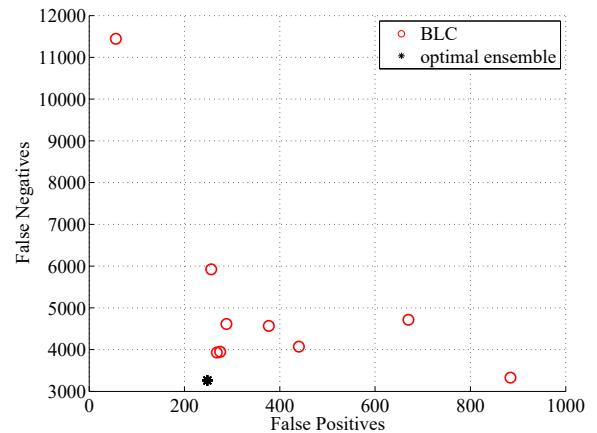


Fig. 2. NSL-KDD Results

B. On UUN Results

As previously discussed, in terms of the UJN dataset, there are only results from the second and third BLC set variations. Across both there are four data points found, each with at least two unique ensembles. The ensembles outperform a number of BLCs in terms of both FP and FN. However, while some of the ensembles outperform a number of BLCs, both PART and JRip produce similar results and one, J48 unpruned, produces the same results.

From Table III and Table IV it may be observed that PART achieved 1 FP with 10 FN to an ensembles 0 FP with FN 10 and JRip achieved 2 FP with 7 FN to an ensembles 1 FP with 6 FN. J48 unpruned achieved the same results as some

TABLE III
INDIVIDUAL BASE LEVEL CLASSIFIER RESULTS

Base Level Classifier	Dataset					
	NSL-KDD		UUJN		Yassin2013	
	FP	FN	FP	FN	FP	FN
NaïveBayesUpdateable	670	4714	67	66	728	37152
PART	275	3948	1	10	11	37159
J48 pruned	263	3900	0	16	60	37159
J48 unpruned	268	3930	3	3	60	37159
DecisionStump	440	4072	65	66	3	37158
DecisionTable	256	5923	92	25	415	37156
ClassificationViaRegression	659	4713				
RandomForest	267	4206				
RandomTree	278	4381				
VFI	377	4567	37	12	54	37069
ConjunctiveRule	884	3328	65	66	0	37159
JRip	288	4613	2	7	11	37159
Nnge	697	3994	430	4		
HyperPipes	56	11443	445	9	211	37069

TABLE IV
OPTIMAL ENSEMBLES AND RESULTS FOUND FOR EACH DATASET AND RELEVANT BLC SET VARIATION

Dataset	Base Level Classifier Set Variation								
	One			Two			Three		
	Chromosome	FP	FN	Chromosome	FP	FN	Chromosome	FP	FN
NSL-KDD	11110000011000	248	3263	01111000000	248	3263	0111101001	248	3263
	11110000000001			01111000001			1111001001		
	01110000010001			01111010000			1111101001		
	11110000010001			11111000001			0111001001		
	11110000000000			01110000000			0111000000		
	01110000010000			11111010000			1111101000		
	11110000010000			01111010001			0111101000		
				11111000000			0111001000		
				0111000001					
UUJN	N/A			11100101011	0	10	0110110001	0	11
				11100101010			0110010001		
				10110101000			1011110001		
				10110101001	1	6	0011110001	1	6
				00110101001			1011110000		
				00110101000			1011010000		
				00111100000			1011010100		
				10011111011	2	4	0011100001	3	3
				00011111011			0001111001		
				10010001011			0001100101		
				00010111000	3	3	0011010001		
				10010101000			0011010001		
				10010101000			1001100101		
				10010011011			0011000000		
				00010101000			1011100100		
				10010011001			1011000100		
				10010111011			1011000000		
				10011101000			1011100001		
						0001000101			
		0001010000							
Yassin2013	N/A						0000100100	0	37159
							1000100100		
							1000100101		
							1000000101		
							0000100101		
							0000000101		
							0000100001	3	37158
							1000100000		
							0000100000		
							1000100001		
							0000111001	43	37069
							0000011001		
							0000011100		
							0000011000		

ensembles of 3 FPs with 3 FNs as evidenced in Tables III and IV. The optimal ensembles achieve low FPs ranging from 0-3. While these values match some FPs produced by BLCs, the FNs indicate that those same instances achieve higher results and, in some cases, the optimal ensembles achieve lower FNs. For example, PART with 1 FP and 10 FNs is outperformed by an ensemble with 0 FPs and 10 FNs. While this difference may be marginal it is still indicative of a better performing ensemble in some cases.

C. On Yassin2013 Results

Finally, with regard to the Yassin2013 dataset, results were only obtained for the third BLC set variation. As there exists only one set of results for this dataset the comparisons between the optimal sets and BLCs are somewhat simplified. Consequently, there are a small number of optimal data points found, each with at least four unique ensembles. Considering the results from each individual BLC, as well as comparisons between the optimal ensembles, it can be observed that some BLCs are able to achieve the same result as optimal ensembles. In particular, this is the case with the BLCs *ConjunctiveRule* and *DecisionStump*, which achieve 0 FPs with 37159 FNs and 3 FPs with 37158 FNs respectively. Comparisons also show that while both the set of individual BLCs and optimal ensembles achieve 37159, 37158 or 37069 FNs, the number of FPs vary more, with the optimal ensembles outperforming most of the BLCs in terms of the number of FPs alone. Optimal ensembles achieve 0, 3 or 43 FPs, where optimal ensembles that achieve 0 or 3 FPs perform better than most individual BLCs where the number of FPs are typically found to be 11 and above, three of which achieve over 200 FPs.

V. DISCUSSION

As stated in Section IV, the optimal ensembles found with NSL-KDD outperform the majority of individual BLCs. Comparing them on their individual FPs and FNs, 100% of the ensembles outperform 93% of the BLCs in regards to FPs while 100% of the ensembles outperform 100% of the BLCs in regards to FNs. Thus, overall, it is shown that 93% of the optimal ensembles strictly outperform any of the BLCs. The greatest distinction between BLCs and optimal ensembles may be the difference found with the single BLC, *HyperPipes*, that is not strictly improved upon; the FPs of the BLC is approximately 22% of any ensemble's FPs, the FNs from the BLC is approximately 350% of any ensemble's FNs. Hence while the FPs may be of more importance for NIDS, the FNs would have more of an effect being instances where an intrusion was not detected. The performances found with UUJN differ in comparison to NSL-KDD, the differences between optimal ensembles and individual BLCs are not as definitive. With four and three possible optimal ensemble results for the second and third BLC set variations respectively, individual comparisons are more varied. For the second BLC set variation at least a single optimal ensemble, for example 11100101011, dominates 63% of individual BLCs. For the third BLC set variation at least a single optimal ensemble,

for example 0110110010, dominates 54% of individual BLCs. The performance when considering Yassin2013 is more easily defined. Two BLCs, *DecisionStump* and *ConjunctiveRule*, achieve results equal to those of two different sets of optimal ensembles. Those optimal ensembles are least likely to dominate given results equal to individual BLCs, thus remains the third optimal ensemble, strictly dominating 60% of the individual BLCs.

Given that adjustments to the sets of BLCs were based on the usability of the datasets, it has been shown that each BLC did not perform equally across each individual dataset. A clear example of this is with the BLC named *ConjunctiveRule*; from results obtained using both the NSL-KDD and UUJN datasets it may be observed that any optimal ensemble outperforms *ConjunctiveRule*. Instead, from results obtained from Yassin2013 it may be observed that *ConjunctiveRule* performs just as well as at least one of the ensembles, for example 0000100100.

Where comparisons across each dataset can be made directly, over the third set variation as given in the last three columns of Table IV, it can be seen that each ensemble found to be optimal is unique; an optimal set of BLC found for NSL-KDD is not found for UUJN or Yassin2013. Accordingly, it may be the case that the data points utilized by each dataset may be less compatible with some BLCs than with others. A particular BLC may perform better when trained on NSL-KDD rather than UUJN; BLCs that do not perform as well may be removed from subsequent generations of chromosomes by the evolutionary process.

When considering the patterns of the presence of the BLCs, over the third set variation, as given in Table IV, it can be seen that for each dataset there is at least one BLC, namely *JRIP*, that does not appear in any of the optimal ensembles. *JRip* is dominated by at least one ensemble whenever used with NSL-KDD, UUJN or Yassin2013. Consequently, it is likely that future experiments may be able to exclude *JRip* as it would appear to play no part in the generation of optimal ensembles.

With the removal of a BLC, *NNge*, between the second and third set variation, there is a slight decrease in performance with UUJN. In the second set variation an optimal ensemble, for example 11100101011, achieves 0 FPs and 10 FNs, while in the third set variation an optimal ensemble, for example 0110110001, achieves 0 FPs and 11 FNs. While these example ensembles differ in the presence of more than one BLC, it may be observed that removal of *NNge* allows production of similar but less favourable results by the remaining BLCs. Furthermore, a different optimal ensemble, for example 1001111011, is no longer produced. Hence, achieving 2 FPs and 4 FNs no longer occurs in the third BLC set variation, perhaps indicating that while the addition of a BLC can improve an optimal ensemble, it may also introduce a wider range of unique optimal ensembles from a set of BLCs.

A comparison between the results obtained and similar works within the literature may also be made. Comparing against [19] shows less FPs however more FNs, all evolved ensembles more accurately detect intrusions but less accu-

rately classify normal instances. It should be noted that the Yassin2013 dataset used in the work presented herein here is only an approximation of the dataset used in [19].

VI. CONCLUSION

The hypothesis was that a GA should be able to evolve sets of BLCs for use with a Multiple Classifier System to produce an ensemble that is able to outperform a number of the individual BLCs from the full set of BLCs. Results show that evolved optimal ensembles can perform better than individual BLCs, the number of optimal ensembles that perform better would appear to partly be based on the dataset.

When considering the BLC set variations, no singular optimal set of BLCs is produced that may be applied with impunity across the varied datasets. However, NSGA-II is capable of evolving sets of BLCs to optimal ensembles for each dataset. Individual BLCs require some analysis before benefits may be produced, especially where prolonged training time requires removals of BLCs, indicating that datasets may dictate inclusion of BLCs. Sets of BLCs may also need to be adequately large to provide optimal variations of optimal ensembles and subsequent data points from evaluation measures. Thus the benefits of the optimization of StackingC could rely on combinations of BLCs and datasets.

Future works could include further investigation and analysis into the selection of individual BLCs better suited for use across varied datasets with hopes that optimal ensembles across varied datasets may be evolved. However, such an approach could be restrictive when attempting to cover a large enough number of datasets. Additionally, implementation of additional GAs and additional objectives could be included as future works.

REFERENCES

- [1] A. Kumar and E. B. Fernandez, "Security Patterns for Intrusion Detection Systems," in *1st LACCEI International Symposium on Software Architecture and Patterns (LACCEI-ISAP-MiniPloP2012)*, Panama City, Panama, 2012.
- [2] M. H. Bhuyan, D. K. Bhattacharyya, and J. K. Kalita, "Network anomaly detection: methods, systems and tools," *Communications Surveys & Tutorials, IEEE*, vol. 16, no. 1, pp. 303–336, 2014.
- [3] M. A. Ambusaidi, X. He, and P. Nanda, "Unsupervised Feature Selection Method for Intrusion Detection System," in *Trustcom/BigDataSE/ISPA, 2015 IEEE*, vol. 1, Aug 2015, pp. 295–301.
- [4] D. Bolzoni, S. Etalle, and P. Hartel, "POSEIDON: a 2-tier anomaly-based network intrusion detection system," in *Fourth IEEE International Workshop on Information Assurance (IWIA'06)*, April 2006, pp. 10 pp.–156.
- [5] P. M. Mafra, V. Moll, J. da Silva Fraga, and A. O. Santin, "Octopus-IIDS: An anomaly based intelligent intrusion detection system," in *Computers and Communications (ISCC), 2010 IEEE Symposium on*, June 2010, pp. 405–410.
- [6] E. Menahem, L. Rokach, and Y. Elovici, "Troika An improved stacking schema for classification tasks," *Information Sciences*, vol. 179, no. 24, pp. 4097 – 4122, 2009.
- [7] A. Ledezma, A. Sanchis, and F. J. Ord, "Genetic Approach for Optimizing Ensembles of Classifiers," in *Proceedings of the Twenty-First International Florida Artificial Intelligence Research Society Conference*, 2008, pp. 89–94.
- [8] M. Tavallae, E. Bagheri, W. Lu, and A. A. Ghorbani, "A detailed analysis of the KDD CUP 99 data set," in *2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications*, July 2009, pp. 1–6.
- [9] A. Shiravi, H. Shiravi, M. Tavallae, and A. A. Ghorbani, "Toward developing a systematic approach to generate benchmark datasets for intrusion detection," *Computers & Security*, vol. 31, no. 3, pp. 357 – 374, 2012.
- [10] Z.-H. Zhou, J. Wu, and W. Tang, "Ensembling neural networks: Many could be better than all," *Artificial Intelligence*, vol. 137, no. 1, pp. 239 – 263, 2002.
- [11] S. Axelsson, "The Base-rate Fallacy and Its Implications for the Difficulty of Intrusion Detection," in *Proceedings of the 6th ACM Conference on Computer and Communications Security*, ser. CCS '99. New York, NY, USA: ACM, 1999, pp. 1–7.
- [12] R. Sommer and V. Paxson, "Outside the Closed World: On Using Machine Learning for Network Intrusion Detection," in *2010 IEEE Symposium on Security and Privacy*, May 2010, pp. 305–316.
- [13] L. Breiman, "Bagging Predictors," *Machine Learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [14] A. Balon-Perin and B. Gambäck, "Ensembles of Decision Trees for Network Intrusion Detection Systems," *International Journal on Advances in Security Volume 6, Number 1 & 2, 2013*, 2013.
- [15] D. H. Wolpert, "Stacked generalization," *Neural Networks*, vol. 5, no. 2, pp. 241 – 259, 1992.
- [16] K. M. Ting and I. H. Witten, "Issues in Stacked Generalization," *J. Artif. Int. Res.*, vol. 10, no. 1, pp. 271–289, May 1999.
- [17] A. K. Seewald, "How to Make Stacking Better and Faster While Also Taking Care of an Unknown Weakness," in *Proceedings of the Nineteenth International Conference on Machine Learning*, ser. ICML '02. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2002, pp. 554–561.
- [18] D. Ariu, R. Tronci, and G. Giacinto, "HMMPayl: An intrusion detection system based on Hidden Markov Models," *Computers & Security*, vol. 30, no. 4, pp. 221 – 241, 2011.
- [19] W. Yassin, N. I. Udzir, and Z. Muda, "Anomaly-Based Intrusion Detection through K-Means Clustering and Naives Bayes Classification," in *Proceedings of the 4th International Conference on Computing and Informatics, ICOCI 2013*, no. 049, 2013, pp. 298–303.
- [20] W. Li, "Using Genetic Algorithm for network intrusion detection," in *In Proceedings of the United States Department of Energy Cyber Security Group 2004 Training Conference*, 2004, pp. 24–27.
- [21] R. H. Gong, M. Zulkernine, and P. Abolmaesumi, "A software implementation of a genetic algorithm based approach to network intrusion detection," in *Sixth International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing and First ACIS International Workshop on Self-Assembling Wireless Network*, May 2005, pp. 246–253.
- [22] 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.
- [23] K. Deb and H. Jain, "An Evolutionary Many-Objective Optimization Algorithm Using Reference-Point-Based Nondominated Sorting Approach, Part I: Solving Problems With Box Constraints," *IEEE Transactions on Evolutionary Computation*, vol. 18, no. 4, pp. 577–601, Aug 2014.
- [24] Q. Zhang and H. Li, "MOEA/D: A Multiobjective Evolutionary Algorithm Based on Decomposition," *IEEE Transactions on Evolutionary Computation*, vol. 11, no. 6, pp. 712–731, Dec 2007.
- [25] M. S. Hoque, M. Mukit, M. Bikas, A. Naser *et al.*, "An implementation of intrusion detection system using genetic algorithm," *arXiv preprint arXiv:1204.1336*, 2012.
- [26] D. S. Kim, H.-N. Nguyen, and J. S. Park, "Genetic algorithm to improve SVM based network intrusion detection system," in *19th International Conference on Advanced Information Networking and Applications (AINA'05) Volume 1 (AINA papers)*, vol. 2, March 2005, pp. 155–158 vol.2.
- [27] G. Kumar and K. Kumar, "Design of an evolutionary approach for intrusion detection," *TheScientificWorldJournal*, vol. 2013, p. 962185, 2013.
- [28] J. J. Durillo and A. J. Nebro, "jmetal: A java framework for multi-objective optimization," *Advances in Engineering Software*, vol. 42, no. 10, pp. 760 – 771, 2011.