

## Methodological Choices in Machine Learning Applications

Kendall E. Nygard  
 Department of Computer Science  
 North Dakota State University  
 Fargo, ND, USA  
 Email: [kendall.nygard@ndsu.edu](mailto:kendall.nygard@ndsu.edu)

Mostofa Ahsan, Aakanksha Rastogi,  
 Rashmi Satyal  
 Department of Computer Science  
 North Dakota State University  
 Fargo, ND, USA  
 Email: {[mostofa.ahsan](mailto:mostofa.ahsan@ndsu.edu), [aakanksha.rastogi](mailto:aakanksha.rastogi@ndsu.edu),  
[rashmi.satyal](mailto:rashmi.satyal@ndsu.edu)}@ndsu.edu

**Abstract**—Machine learning is a subset of artificial intelligence in which a machine has an ability to learn and employ complex algorithms to impersonate human behavior. Development of a machine learning model involves careful preparation and management of data and selection and features to produce meaningful results. The data issues are often challenging due to availability, characteristics, properties, categorization, and balance. We report on relevant literature, case studies and experiments surrounding the data issues. We describe alternative machine learning methodologies and emphasize supervised learning, including treatment of experimental procedures. Procedures and challenges in the collection, quantity, distribution, quality, sampling, of and relevancy of data are included. Applications of machine learning models are presented, including classification models for self-driving cars. These models introduce anti-autonomy trust modeling. We also describe intrusion detection models that can detect malicious activity in computing systems. These applications also provide insight into overfitting and underfitting training data. Feature engineering and feature selection issues are presented, including approaches to identifying, combining, and eliminating attributes and features to determine which are needed and their significance. Approaches for treating class imbalances in data management are discussed. Comparisons among categorical encoding techniques are presented. The work provides perspective and insight into resolving multiple issues that must be addressed in utilizing machine learning models in practice.

**Keywords**- *Machine Learning; Data Management; Feature Engineering; Feature Selection; Self-Driving Cars; Intrusion Detection.*

### I. INTRODUCTION

Machine Learning (ML) is a rapidly emerging area of artificial intelligence. Many types of applications have been successfully developed and new successes are regularly reported. The famous Turing award winner Jim Gray referred to data science as a “fourth paradigm,” taking a rightful place among empirical, theoretical, and computational sciences [2].

Often viewed as interdisciplinary, data science involves mathematics, statistics, and computer science as well as other related areas. In many applications, the availability of large and relevant datasets, and the methods of data science provide the lifeblood of machine learning problem-solving approaches. Analyses and decision support in nearly every area of human endeavor today are related to machine learning.

The example machine learning studies that we describe are in the areas of self-driving cars and intrusion detection [1][3][34][37][38]. Self-driving cars are a real-world example of a system that requires machine learning, since they involve complex computations, algorithms, computing systems, mechanics, and behavioral aspects that endanger human life if automated decisions or controls go awry. Automated cars have features such as remote engine start, advanced information systems, moving object detection (MOD), lane change assist, anti-lock braking system (ABS), to name just a few. These systems can create vulnerabilities to cyber-attacks from things like bugs introduced in their core software code, remote access to the onboard diagnostic system (OBD) of the vehicle, or controller area network (CAN) bus [77].

For a self-driving car study focused on classification, we utilized available multi-attribute data about specific collisions. The data contained many features and attributes of the vehicle itself, the damage incurred, roadway conditions, etc. The objective of the study was to build a classification model that could translate the detailed data into collision predictions and to drive an anti-autonomy trust model. There were several important and difficult choices made related to scale and balance within the available dataset, and in feature engineering. A linear sequential supervised machine learning model was employed.

The intrusion detection study used supervised learning techniques to build a model for identifying outside threats initiated by malicious actors who wish to breach or compromise a system. Among other datasets, the study examined the famous dataset that originated in the KDD (Knowledge Discovery and Data Mining) competition and was later modified to form the now publicly available NSL (Network Security Laboratory) KDD dataset [7][8].

Data management and feature engineering are important steps in the cybersecurity domain to prepare data for machine learning algorithms and build effective models for detecting security threats and anomalies. Data analysis involves exploring, visualizing, and understanding the data to identify patterns, trends, and anomalies [48]. This may include statistical analysis, correlation analysis, time series analysis, and other techniques to gain insights into the data. In the cybersecurity domain, data analysis may also involve extracting features that are relevant to the security domain, such as network traffic flow, packet size and content, system

logs, user behavior, and other indicators of security threats [49]. Feature engineering is the process of selecting, transforming, and creating features that are relevant and informative for the machine learning model. This may involve selecting features that are correlated with the target variable, transforming features to make them more informative or meaningful, and to create new features from existing ones [50]. Feature engineering in the cybersecurity domain may also involve domain-specific knowledge and expertise to select and create features that capture the specific characteristics of security threats and anomalies [51][52]. Feature engineering techniques may also involve dimensionality reduction, feature scaling, and other preprocessing techniques to improve the performance of the machine learning model. All these methods have their strengths and weaknesses, and their effectiveness may depend on the specific dataset and classification problem at hand. In the cybersecurity domain, they can be useful to improve the performance of machine learning models that aim to detect and classify security threats, anomalies, or attacks [53]. We describe some of the data analysis and feature engineering techniques used in cybersecurity in the data management section of this paper.

The rest of the paper is structured as follows. In Section II, we describe supervised machine learning with illustrations of the flow of a machine learning model and data splits for cross validation. In Section III, we present a self-driving cars example illustrating an implementation of a linear sequential supervised learning artificial neural network model utilizing multiple pre-processed complex attributes. In Section IV, we present the intrusion detection example, explaining how a machine learning model can be tuned to predict and identify attacks. In Section V, we describe data management and feature engineering issues that are ubiquitous in machine learning practice. This section also includes several categorical encoding techniques for preprocessing data for a machine learning algorithm. Finally, we conclude our work in Section VI. An abbreviated version of the work is available in [1].

## II. MACHINE LEARNING EXPERIMENTS

Machine learning methods are of four distinct types. Supervised learning models are trained using datasets with known labels, then used to make predictions on new data. Unsupervised learning models work with unlabeled data and seek clusters or patterns in the data. Semi-supervised learning is a hybrid approach that uses both unsupervised and supervised learning techniques. Reinforcement learning uses no labeled data, but instead is based upon evaluation of rewards or punishments of behaviors. We have conducted extensive experiments using supervised learning in applications concerning circumstances under which collisions occur in self-driving cars and in detecting intrusions into computer platforms.

Supervised machine learning (SML) methods are very effective in addressing classification problems. When applied to classification tasks, a SML method has a set of available

data for which their correct classifications are known. Such a data set can be represented as shown below.

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots (\mathbf{x}_n, y_n)$$

Here, the  $\mathbf{x}_i$  shown in bold are vectors that capture a data instance or situation, and the corresponding  $y_i$  values are discrete labels for the available classifications.

The initial task in a supervised learning experiment is to computationally train the machine learning model to accept the known data instances as input and to produce the correct target as output. Many types of SML methods can be used in training, including decision trees, neural networks, support vector machines, and logistic regressions [78]. For some methods, training can be a computationally intensive process. Once trained, the model is available to accept new data instances and predict their target classification. The model is successful if it has high values of performance measures such as percentage of accuracy in correctly classifying the new data instances, called the ability to generalize. There are multiple issues surrounding the characteristics of the available data, the classes into which they fall, their attributes and features, and the learning models charged with producing the predictions. Concerning baseball, Coach Yogi Berra famously said, "It's tough to make predictions, especially about the future." This aphorism is equally true in machine learning [46].

Figure 1 illustrates the general flow of a machine learning technique. Several tasks are included. The overall task of the DEVELOP phase shown at the top is to produce a Final Model that is fully specified, trained and feeds into the PREDICT phase shown below the dotted line, where it is available for generalization use on new data. Starting from the top, the data is shown as partitioned into splits for Training, Validation and Test. The full data is divided into the Training Split and the Test Split. A good way to perform training is to withhold a portion of the data while training is done. It is viewed as a mistake to train and test a machine learning model on the same data. So, doing that would result in the model memorizing all the data/target pairs, resulting in the model perfectly knowing all of the answers, leaving no ability to generalize. The result is known as *overfitting*. For validation purposes, the Training Split is typically divided into pieces called folds. Called k-fold cross validation, Figure 2 illustrates basic logic for splitting the data. In this example,  $k=5$  so there are five equal parts. This corresponds to the Validation Split and Model Tuning blocks in Figure 1. In Figure 2, shown in bold italics on the diagonal, there is a designated fold in each row that is specified for testing, with the other four used for training. The key idea is to find the best set of meta-level parameters for a model being developed. All major machine learning models have parameters. For example, an ML that utilizes an Artificial Neural Network (ANN) in some way, will be parameterized with settings like Learning Rate (governs weight adjustments), topology (number of hidden layers, nodes

within layers, and interconnectedness), and activation functions.

In cross validation, when a model is trained on the folds, a performance metric, such as classification accuracy, can be calculated on the testing fold. After all the fold splits are evaluated in this way, an average is calculated, which yields a score for the parameter settings. Various optimization methods can be employed to explore the parameter space in a quest to identify the best settings. Viewed more generally, the Model Tuning block can also be viewed as exploring various types of models in a quest to not only optimize the use of one type of model, but to also choose among competing models.

In multiple places of the ML process, there is a need to evaluate the quality of the predictions using a metric. The empirical accuracy of a method is simply the percentage of the predictions made that are correct. Other metrics are available. More details are provided later in this paper.

Raw data is rarely available in a form that is suitable for direct use by an ML model. Pre-processing of data is typically necessary to deal with such things as null or missing values, outliers, transforming or reconciling numeric and categorical values, rescaling, and standardizing. We expand on the data-centric issues, for the example applications reported in this article.

Feature Engineering also appears in Figure 1. Features are those characteristics that are present in the data that are potentially useful in predicting a target outcome.

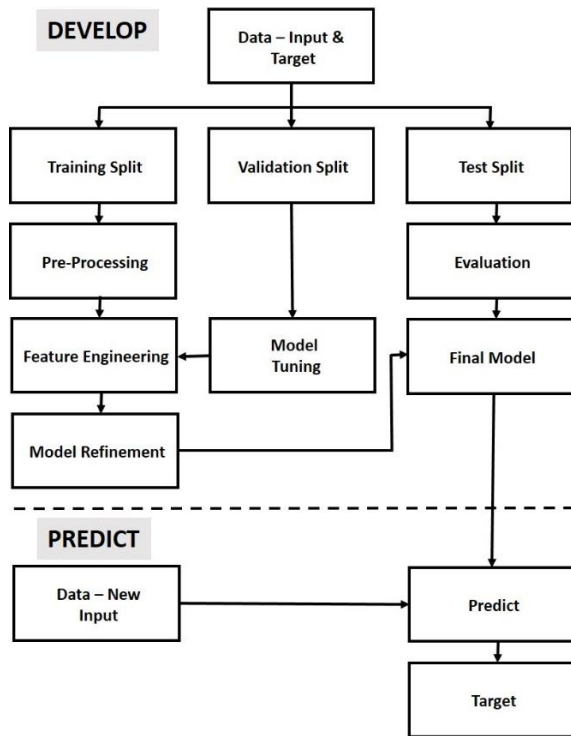


Figure 1. Flow of a Machine Learning Model [3].

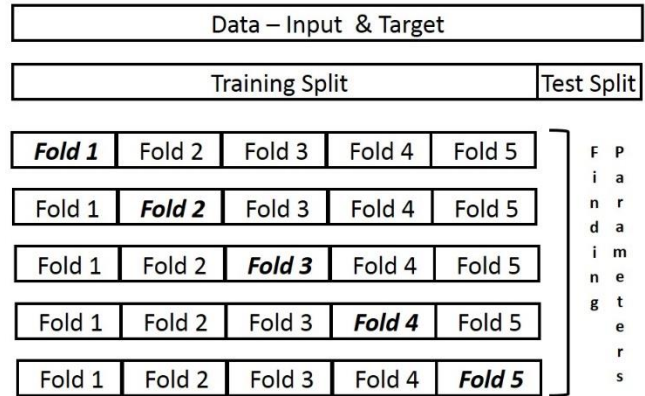


Figure 2. Data Splits for Cross Validation.

It is often effective in ML to modify or combine features in some way to produce a new feature that can improve the prediction accuracy of the method. Called Feature Engineering, the operations that can be carried out include things like mathematically transforming a single feature or applying a functional calculation on multiple features. Feature Selection refers to reducing the number of features employed by the model while retaining acceptable results. Reducing the features needed can ease the data collection task and reduce the computational load of running the model. Feature Selection typically follows Feature Engineering. We provide details related to the examples discussed in this paper.

Unsupervised Learning is different than supervised learning in many ways. Some of the most known algorithms are, k-means clustering, hierarchical clustering, principal component analysis, and a priori algorithm [45]. The need for unsupervised models is increasing in the cybersecurity domain since attacks are being modified every day [47].

We have discussed multiple machine learning techniques in this section. The primary concern is to make proper choice of methods to optimize the solution of a problem. We discuss the criteria we need to adopt to address a machine learning problem in the following sections.

### III. SELF-DRIVING CARS APPLICATIONS

A self-driving car underway must adhere to laws and rules of the road, like adhering to speed limits, stopping at red traffic lights, turning from an appropriate lane, etc. In addition, the vehicle must carry out control actions in response to sensor and communication information that provides information regarding things like conditions of the roadway (such as snow, rain, deteriorated pavement); construction zones; presence of bicycles, pedestrians, other cars, or obstacles; or visibility issues like glare, fog, snow, rain, darkness, or any type of impaired lighting.

When addressing how a self-driving car can be trained to carry out an appropriate action under circumstances that it encounters, we draw an analogy with how positive

reinforcement can work efficiently when a person or animal learns a new skill or behavior. This vehicle training need lends itself naturally to Reinforcement Learning (RL), a powerful machine learning technique that rewards good behaviors, and as needed, punishes bad behaviors. A RL training method proceeds iteratively and is illustrated in Figure 3 [76]. While undergoing learning in a simulated environment, the vehicle is an agent that carries out actions, receives their consequences in the form of positive or negative feedback, and adjusts their model of actions accordingly. When complete, experience gained in this way can maximize rewards. With acceleration, deceleration and steering as the primary actions, a self-driving car (agent) aims at maximizing short-term rewards (like safe driving) and long-term rewards (like fast arrival at the destination) using the RL approach.

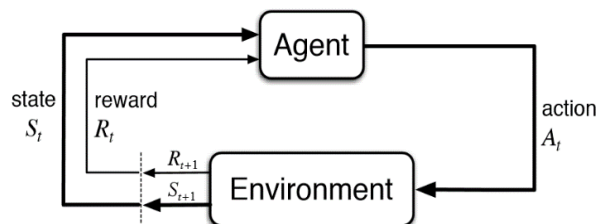


Figure 3. Reinforcement Learning illustration.

We carried out extensive experiments using supervised learning for analyses of collisions that occur with self-driving cars. Official collision reports basically map data items that describe driving conditions into a collision classification. We used these reports to help determine circumstances under which self-driving cars carry out actions that cause collisions or fail to avoid them. These harmful actions are referred to as anti-autonomy traits and factors on the part of the vehicle that cause collisions and diminish trust [3][39]. Data availability was a challenge since jurisdictions of different states, federal traffic agencies and motor vehicle departments often do not make their data publicly available. Data used in this study was submitted by the manufacturers of autonomous vehicles to the California Department of Motor Vehicles for collisions that occurred with other cars, pedestrians, bicyclists, and other objects during their test drives on roads and freeways in the state. All data applied to collisions that occurred while the cars were being driven in autonomous driving mode. The collisions occurred between October 2014 and March 2020 [3][39]. The attributes of this dataset are listed in Table I below. All attribute names are feature type categorical and data type object.

TABLE I. COLLISION DATA ATTRIBUTES [3]

Attribute Type	Attribute Names
Autonomous vehicle details	Manufacturer Name, Business Name, Vehicle Year, Vehicle Make, Vehicle Model, Vehicle was (stopped in traffic/moving)

Attribute Type	Attribute Names
Accident Details	Date of Accident, Time of Accident
Involved in Autonomous vehicle accident	Involved in Autonomous Vehicle Accident (Pedestrian/Bicyclist/Other), Number of vehicles involved with Autonomous Vehicle
Autonomous vehicle damage	Vehicle Damage, Damaged Area
Details of other vehicle involved in accident	Vehicle 2 Year, Vehicle 2 Make, Vehicle 2 Model, Vehicle 2 was (stopped in traffic/moving)
Involved in other vehicle accident	Involved in Vehicle 2 Accident Pedestrian, Involved in Vehicle 2 Accident Bicyclist, Involved in Vehicle 2 Accident Other, Number of vehicles involved with Vehicle 2
Injuries	Injured, Injured Driver, Injured Passenger, Injured Bicyclist
Vehicle driving mode	Vehicle Driving Mode
Weather conditions for both vehicles	Clear, Cloudy, Raining, Snowing, Fog/Visibility, Other, Wind
Lighting conditions for both vehicles	Daylight, Dusk-Dawn, Dark Street Lights, Dark-No Street Lights, Dark-Street Lights Not Functioning
Roadway surface for both vehicles	Dry, Wet, Snowy-Icy, Slippery/Muddy/Oily/etc., Holes-Deep-Rut, Loose Material on Roadway, Obstruction on Roadway, Construction/Repair Zone, Reduced Roadway Width, Flooded, Other, No Unusual Conditions
Preceding Movement of Autonomous Vehicle before collision	Stopped, Proceeding Straight, Ran Off Road, Making Right Turn, Making Left Turn, Making U Turn, Backing, Slowing/Stopping, Passing Other Vehicle, Changing Lanes, Parking Maneuver, Entering Traffic, Unsafe Turning, Xing into Opposing Lane, Parked, Merging, Travelling Wrong Way, Other
Preceding Movement of Other Vehicle before collision	Stopped, Proceeding Straight, Ran Off Road, Making Right Turn, Making Left Turn, Making U Turn, Backing, Slowing/Stopping, Passing Other Vehicle, Changing Lanes, Parking Maneuver, Entering Traffic, Unsafe Turning, Xing Into Opposing Lane, Parked, Merging, Travelling Wrong Way, Other
Type of Collision	Head On, Side Swipe, Rear End, Broadside, Hit Object, Overturned, Vehicle/Pedestrian, Other
Other	CVC Sections Violated Cited, Vision Obscurement, Inattention, Stop and Go Traffic, Entering/Leaving Ramp, Previous Collision, Unfamiliar with Road, Defective WEH Equip Cited, Uninvolved Vehicle, Other, Non-Apparent, Runaway Vehicle

This data was extracted from PDF files and converted into CSV format with 256 rows of data in 140 columns. As is often the case in machine learning, data cleaning was a significant effort, and included pre-processing steps for augmenting, labeling, and classifying the data [3][39]. Data was augmented to 5256 rows with the goal of yielding a model with minimal noise. Details of data augmentation criteria is described in [3]. The core purpose of the study was to associate conditions into a level of trust that people had in a self-driving car. The values of attributes in the data that describe conditions and circumstances that are present when a collision occurs provide a handle to model a mapping between data and trust level. After pre-processing the data, a linear sequential supervised learning ANN model called NoTrust was devised, validated, and tested to classify the target data, using the basic approach illustrated in Figure 1.

The model used the libraries provided by Keras with the Tensorflow backend [40][41][42]. Python was used for programming since it integrates with Keras to access the neural network Application Programming Interface. The deep learning libraries of Keras support fast prototyping, modularity, and smooth computation.

There are multiple challenges concerning data, features, and metrics in applying the ML methodology to the self-driving car application. First, there were only 256 collision reports available, which is arguably a small number to use in a ML method. This was mitigated by augmenting the data to produce a larger set to develop a model with higher accuracy. Also, in the context of alternative target value possibilities, the data is unbalanced in that the number of samples across the distinct classes differs widely. Section V describes methods, such as oversampling, to deal with unbalanced data. Second, there are 140 attributes, a large number relative to sample size, as shown in Table I. Thus, the possible permutations and combinations that could be evaluated in the ANN model are explosive. Fortunately, with initial analyses of the data and evaluation runs, it was possible to identify a subset of attributes and features that were mandatory to include. This analysis was done by systematically configuring sets with and without specific attributes, evaluating each combination, and comparing the outcomes. Using this search method, we identified certain categorical features that were closely correlated with the anti-autonomous traits of the vehicles. This supported dropping the features related to date/time of accident, vehicle manufacturer, weather conditions, lighting conditions, roadway conditions, vehicle movement and type of collision. Ultimately, we arrived at the five attributes shown below to form the mandatory set.

- Vehicle driving mode = autonomous
- Vehicle damage = moderate and major
- First vehicle involved = Pedestrians/Bicycle/Other
- Second vehicle involved = Bicycle/Other
- Injuries sustained = Pedestrians/Bicyclists/Others

While keeping the model simple and still retaining accuracy, the mandatory feature set performed well in making trust and do not trust predictions for autonomous vehicles. However, when anti-autonomous traits of the self-driving car itself were incorporated into the model it became apparent that more attributes had to be utilized.

Anti-autonomy refers to decisions and actions taken by a self-driving car that are in some way inappropriate in terms of increasing risk, diminishing safety, causing potential harm, or lowering trust. It entails from an unexpected, unconventional, and abnormal decisions that self-driving car makes in the event of unfavorable surrounding driving conditions related to weather, roadway surface conditions, drivability of other vehicles and, pedestrians or bicyclists sharing the road. Autonomous vehicles have been known to have certain technological shortcomings in terms of Lidar failing when the weather conditions are rainy or foggy with limited visibility. Extensive study and analysis of the collision data revealed the following anti-autonomous behavior of self-driving cars [3]:

- In 50.22% of collisions, a pedestrian was involved while the vehicle was driving autonomously.
- In 52.08% of collisions, a pedestrian was involved while vehicle was driving autonomously and was moving in traffic.
- In 55.13% of the collisions, a bicyclist was involved while the vehicle was driving autonomously and was about to slow down or stop. These statistics illustrate the confusing behavior of a self-driving car, its mechanical functioning and decision making at the second when a bicyclist appears in front of them.
- In 55.84% of the collisions, a pedestrian was involved while the vehicle was driving autonomously and was attempting to make a parking maneuver. These statistics reveal a potential malfunction of vehicle operation in terms of gauging pedestrian behavior and hitting the breaks just in time.
- 54.28% of the collisions happened when the vehicle was driving autonomously during foggy weather with limited visibility. These statistics depict the malfunctioning of the sensors, camera, and Lidar sensors.
- In 51.71% of the collisions, a pedestrian was involved while the vehicle was driving autonomously during foggy weather with limited visibility.
- In 50.68% of the collisions, a bicyclist was involved when the vehicle was driving autonomously at night when streetlights were on.
- In 53.18% of the collisions, a pedestrian was involved when the vehicle was driving autonomously at night when streetlights were on.

The anti-autonomous traits that were incorporated into the model include vehicle driving mode, type of collision, weather conditions, roadway surface conditions and injuries sustained by pedestrian/bicyclists/others. In addition to the linear sequential ANN, evaluation of Recurrent Neural Networks (RNN) models with Long Short-Term Memory (LSTM) were available for possible comparison purposes. The primary reason for choosing a sequential ANN model was that the classification sought is binary, predicting whether the autonomous vehicles could be trusted. A sequential ANN model utilizes a stack of layers with each layer having exactly one input tensor and one output tensor. This contrasts with a Functional API with shared layers, non-linear topology, and multiple inputs/outputs. This study uses an input stack of layers of selective features that have a single output to model affirmation or denial of trust. Thus, we avoided layer sharing, non-linear topologies, multiple inputs/outputs, or creation of a directed acyclic graph or graph of layers. These properties favor the choice of sequential ANN. Also, this model was chosen over RNN because the data utilized for model processing was not a time series or natural language sequence data. When the additional attributes are included, along with measures of severity of damage sustained by vehicle, the imbalance of the data increases. More specifically, the larger number of predictors added more noise, redundancies, increased overfitting, and decreased the quality of the predictions. A related study by Meiri and Zahavi [4] used simulated annealing to search the attribute space.

Combinatorial problems often have issues related to model accuracy, performance, and optimizer bias. Also, the model solutions offered by machine learning include approximation errors, which further exacerbates issues related to differences between training and validation data [5]. This can be solved by two approaches – active learning and passive learning. Active learning involves updating the model itself to assure a convergence between training and validation curves, in turn improving model accuracy and optimization bias. Passive learning involves the training set providing a uniform and sufficient coverage of the search space [5]. In a similar context, Charikar, Guruswami, Kumar, Rajagopalan, and Sahai [6] defined and studied combinatorial feature selection problems, presented a theoretical framework, and provided algorithms on approximation and hardness results of these combinatorial problems [6].

#### IV. INTRUSION DETECTION APPLICATIONS

In today's world of connected devices, security of the network is of critical importance. Unauthorized access and malicious activities are a great threat to confidentiality, integrity, and availability that form the information security triad. The role of an Intrusion Detection System (IDS) is to detect abnormalities caused by an unauthorized reach into the network and send alerts. An IDS is an element of support for a wall of defense between cyber-attacks. Supervised ML

techniques in an IDS can provide high detection accuracy, particularly against known types of attacks.

The NSL-KDD is an update and improvement to the KDD'99 dataset that was developed for the KDD Cup competition in 1999 [7]. These datasets are publicly available and are very widely used for IDS experiments. The data is primarily internet traffic consisting of 43 features per record, of which the last two are *class* (attack or normal) and *score* (severity of traffic input) [8]. The *class* column provides information on whether the record is considered normal or is a member of one of four attack classes - Denial of Service (DoS), Probe, Remote-to-Local (R2L) or User-to-Root (U2R). There are 14 attack types under these 4 classes: Apache2, Smurf, Neptune, Back, Teardrop, Pod, Land, Mailbomb, Processtable, UDPstorm, WarezClient, Guess\_Password, WarezMaster, Imap, Ftp\_write, Named, Multihop, Phf, Spy, Sendmail, SmpGetAttack, AnmpGuess, Worm, Xsnoop, Xlock, Buffer\_Overflow, Httptuned, Rootkit, LoadModule, Perl, Xterm, Ps, SQLattack, Satan, Saint, Ipsweep, Portsweep, Nmap, Mscan [43][44]. A mixture of categorical (nominal), binary and numeric variables are in the feature set. Each record has basic, content-related, time-related, and host-based features [9]. The attributes of this dataset are listed in Table II.

TABLE II. NSL-KDD DATASET ATTRIBUTES [9]

Attribute Type	Attribute Names
<b>Basic</b>	Duration, Protocol_type, Service, Flag, Src_bytes, Dst_bytes, Land, Wrong_fragment, Urgent
<b>Content related</b>	Hot, Num_failed_logins, Logged_in, Num_compromised, Root_shell, Su_attempted, Num_root, Num_file_creations, Num_shells, Num_access_files, Num_outbound_cmds, Is_hot_login, Is_guest_login
<b>Time related</b>	Count, Srv_count, Error_rate, Srv_error_rate, Rerror_rate, Srv_error_rate, Same_srv_rate, Diff_srv_rate, Srv_diff_host_rate
<b>Host based traffic</b>	Dst_host_count, Dst_host_srv_count, Dst_host_same_srv_rate, Dst_host_diff_srv_rate, Dst_host_same_src_port_rate, Dst_host_srv_diff_host_rate, Dst_host_error_rate, Dst_host_srv_error_rate, Dst_host_error_rate, Dst_host_srv_error_rate

The study also used the UNSW-NB15 dataset. This dataset has 49 features categorized into 6 groups: basic, flow, time, content, labelled and additional generated features [10]. There are 9 attack types: fuzzers, analysis, back-doors, DoS, exploits, generic, reconnaissance, shell code and worms [11]. This dataset has a mixture of categorical, binary, and numerical datatypes. The attributes of this dataset are listed in Table III below.

TABLE III. UNSW-NB15 DATASET ATTRIBUTES [16]

Attribute Type	Attribute Names
<b>Basic</b>	state, dur, sbytes, dbytes, stl, dtl, sloss, dloss, service, sload, dload, spkts, dpkts
<b>Flow</b>	srcip, sport, dstip, dsport, proto

<b>Content</b>	swin, dwin, stcpb, dtcpb, smeansz, dmeansz, trans_depth, res_bdy_len
<b>Time</b>	sjit, djit, stime, ltime, sintpkt, dintpkt, tcprtt, synack, ackdat
<b>Additional generated (general purpose)</b>	is_sm_ips_ports, ct_state_ttl, ct_flw_http_mthd, is_ftp_login, ct_ftp_cmd
<b>Additional generated (connection)</b>	ct_srv_src, ct_srv_dst, ct_dst_ltm, ct_src_ltm, ct_src_dport_ltm, ct_dst_sport_ltm, ct_dst_src_ltm
<b>Labelled</b>	attack_cat, attack_cat

The target attribute either identifies records as ‘normal’ or ‘attack’ or distinguishes the record as a particular attack type. Depending on the desired goal of an intrusion detection system, the machine learning model is tuned to identify a particular attack, which is a challenge. It is thus essential to understand the requirement thoroughly and preprocess input data accordingly.

As an illustration of evaluation metrics, at a high level in the IDS study, for each input vector we have exactly one of the following outcomes:

TP = True Positive = Correct predication that the input vector is an Attack

TN = True Negative = Correct prediction that the input vector is not an Attack

FN = False Negative = Incorrect prediction the input vector is not an Attack

FP = False Positive = Incorrect prediction the input vector is an Attack

The most widely reported metric is Basic Accuracy of the model, which simply reports the proportion of attack reports that are correct.

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

Basic Accuracy is notoriously deceptive when the classes are unbalanced, as in the case of intrusion detection studies, where most input vectors are not attacks. False reports are of interest. This gives rise to the need for metrics such as Precision and Recall, which can be calculated from information in the confusion matrix given below.

	Prediction		
	Attack	Not an Attack	
Actual	Attack	TP	FN
	Not an Attack	FP	TN

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

Precision measures the proportion of the vectors reported by the IDS as attacks that are real attacks. Recall measures the proportion of the vectors that are real attacks and do get reported as such. This means that when Recall is high the IDS does not misclassify many true attacks.

In Intrusion Detection applications, false negatives can be very deadly, which favors high Recall. However, dealing with false positives also has a cost. Unfortunately, experiments that improve Precision typically reduce Recall. The reverse is also often true. For this reason, the harmonic mean of the two, called the F1 score is often calculated.

$$F1 = (2 * (\text{Precision} * \text{Recall})) / (\text{Precision} + \text{Recall}).$$

The effect of the F1 score, which falls between 0 and 1, is to punish extreme values.

The NSL-KDD and UNSW NB-15 datasets are used for training machine learning models by several researchers. In [34] researchers trained the KDD-99 dataset with a mutation of a convolutional neural network and Long Short Term Memory (LSTM) network. The machine learning algorithm test accuracy was 99.7% , which outperformed other models, including DenseNet, CNN (Convolutional Neural Network), GRU (Gated Recurrent Unit), BiLSTM, and Auto Encoder. The experiment was multilabel and nearly all the individual target variables had f1-score, precision and recall exceeding 98%. Work reported in [54] and [55] took similar steps for setting up a machine learning model experimental design to train the NSL KDD and UNSW NB-15 datasets and were able to improve accuracy in both cases.

The NSL-KDD data set has 49 attributes of 6 types and 9 types of intrusions. We considered possibilities for reducing the number of attributes without eliminating information critical to the classifications. Principal Component Analysis (PCA) is a time-honored statistical method for identifying high correlations and ranking the relative importance of the attributes. Although we considered PCA, we alternatively chose to implement an autoencoder neural network. In brief, an autoencoder is a multiple layer neural network in which the input and output layers are identical, and a middle layer is of smaller dimensionality. An autoencoder is a deep learning unsupervised learning method in that the labels for the known data play no role, and after training, the middle layer becomes a compressed version of the input data. The middle layer serves as a pattern that is a discrete model of the data in a compact form. The approach requires experimentation that searches through alternative topologies and tuning parameters of the neural network, including number of layers, nodes within layers, learning rate, and number of epochs. Results of the autoencoder experimental work are given in Table IV below.

TABLE IV. PERFORMANCE EVALUATION OF IDS CLASSIFICATION

Algorithm	Accuracy (%)	Precisio n	Recall	F1_score
Autoencoder	88.76	0.852	0.971	0.908
SVM	86.54	0.824	0.913	0.901
Logistic Regression	82.12	0.808	0.921	0.893

The results show that the model performs well as a binary classifier for threat detection. After applying auto encoder, the accuracy reaches 88.76%, False positives occur at a low rate, so that the method flags only a few normal network inputs as attacks, which is useful in keeping a focus on real threats. Importantly, recall at a high level of 97.1% means that nearly all real attacks are correctly identified and flagged. This is a must-have feature of an IDS since undetected attacks can be very damaging. SVM with RBF kernel (Radial Basis Function) competes with the Autoencoder with a good F1 score. The simple logistic regression works well for binary classification but fails to achieve better results for the NSL KDD dataset. These experiments are a baseline for machine learning methodologies and their proper application. We did not perform any parameter tuning for this experiment, which might have improved the accuracy and performance metrics.

## V. DATA AND FEATURE ENGINEERING

We provide descriptions of data management and feature engineering issues that are pervasive in ML practice and were of importance in our applied studies.

### A. Data Management

Class imbalance in a dataset means that the relative numbers of instances within the classes vary significantly in number [17]. The magnitude of the discrepancies will also vary. Class imbalance is common in most important data domains, including detection of things like fraudulent activities, anomalies, oil spills, and in medical diagnoses. The imbalance of classes occurs in both binary class and multi-class classification [18]. In binary classes, the smaller and larger cardinality classes are called minority and majority classes respectively [17][19]. Class imbalance can influence the training process of ML techniques and lead to ignoring the minority class entirely. We discuss some of the approaches to treat class imbalances. Figures 4 – 10 illustrate the results of applying each technique.

**Random oversampling of a minority class.** In this approach data instances in minority classes are duplicated at random to induce a balance of membership between classes. Due to randomness of the oversampling, the method is naïve in that it makes no assumptions about the classes and their members [20][21]. Since exact copies of some data instances are included in training, there is a risk of overfitting. Classifier accuracy may also be influenced, and computational effort may be increased.

**Random undersampling of a majority class.** This approach discards data instances from majority classes to induce balancing [22]. As in the case of the oversampling method, the discarded data is chosen randomly and naively. The method applies to both binary and multiclass data. The approach can make it difficult to distinguish boundaries

between classes, with an inimical impact on performance measures [23].

**Synthetic Minority Oversampling Technique (SMOTE).** This technique was introduced in 2002 to address the shortcomings of the oversampling and undersampling approaches [24][25]. The technique generates synthetic data by calculating feature space similarities between minority class data instances. The K-nearest neighbors of each data instance in a minority class are calculated, then randomly selected one by one. The method then calculates linear interpolations among the data and uses them to create synthetic data instances.

**Borderline SMOTE.** The SMOTE approach encounters issues when minority class data instances occur in the vicinity of majority class data instances, creating undesirable bridges. The Borderline SMOTE variation addresses this drawback by oversampling only minority class instances near the borderline. Data points are called border points if they are incident to both minority and majority classes and called noise points otherwise [26]. Border points are then utilized to balance the data between classes.

**K-Means SMOTE.** This technique generates minority class samples in safe and crucial borders of input spaces and thus assists performance in classification. The method begins by clustering the dataset using the K-means procedure, then selects the clusters that have higher numbers of minority samples [27]. Additional synthetic samples are then assigned to clusters where minority class samples are sparsely distributed. No noise points are generated.

**SVM SMOTE.** A variation of Borderline-SMOTE, the method finds misclassification points. The borderline points are approximated and classified with a Support Vector Machine (SVM) classifier [28]. Synthetic data points are created randomly around the lines joining each minority class support vector with its neighbors.

**Adaptive Synthetic Sampling – ADASYN.** A limitation of Borderline SMOTE is that it utilizes only synthetic points generated from extreme observations and the borderline instances and neglects the rest of the points in minority classes. This issue is addressed by ADASYN by creating synthetic data using the density of the existing data [29]. The ratio of synthetically generated data is created in inverse relation to the minority class density. In this way, a less dense area creates more synthetic data.

The Churn Modeling Data from Kaggle was applied to the methods [30]. Figure 4 shows the distribution of the data in the original classes, followed by the outcomes of the alternative methods in Figures 5 to 10.



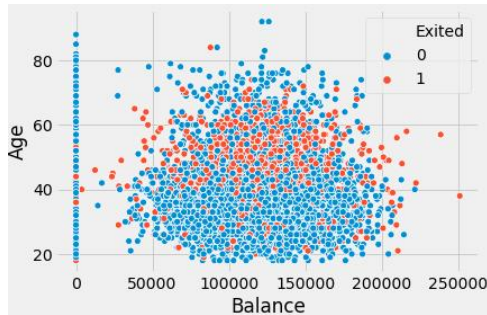


Figure 4. Original Class Imbalance Illustration.

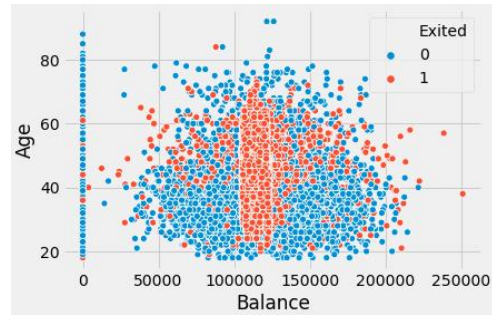


Figure 8. Outcome of K-Means SMOTE.

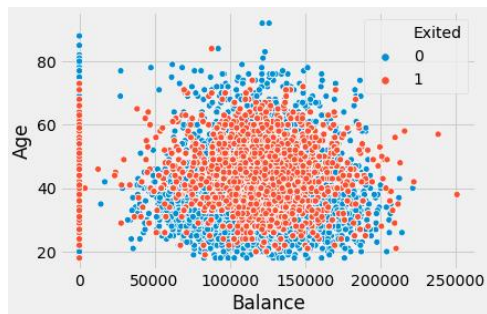


Figure 5. Outcome of Random Oversampling.

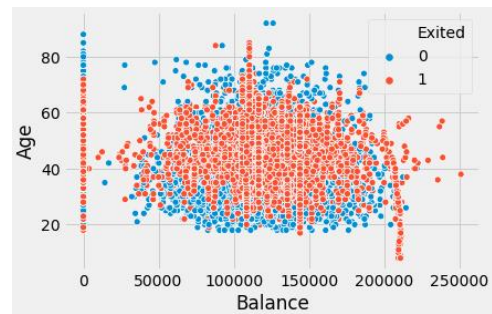


Figure 9. Outcome of SVM SMOTE.

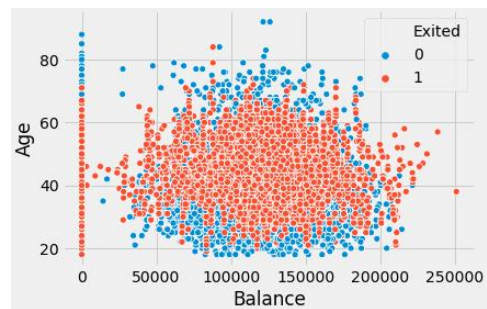


Figure 6. Outcome of SMOTE.

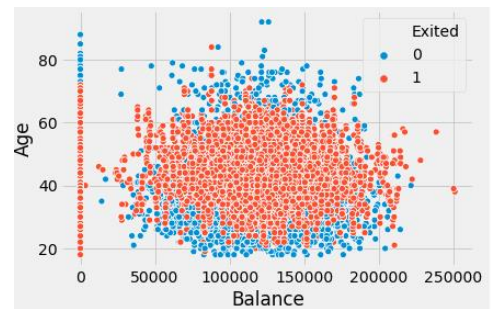


Figure 10. Outcome of ADASYN.

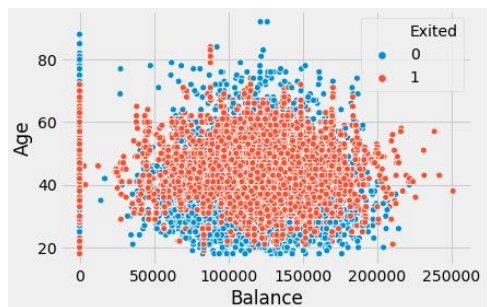


Figure 7. Outcome of Borderline-SMOTE.

### B. Data Management and Feature Engineering

There are multiple methods for feature engineering on categorical data. The inputs to ML algorithms must be numeric, but many applications have categorical data. In our work with ML methods for self-driving car collisions there are examples of ordinal categorical data, such as rating of severity of damages or a weather condition. The intrusion detection work examples include counts of file access attempts, session duration, or error rates. There are also examples of nominal data, such as a type of vehicle in our self-driving car work, or whether a flag is set in the intrusion detection work. Various encoding methods are used to convert the variables into a useful numerical representation [10]. Choosing an appropriate encoding scheme is an essential part of data preprocessing for a ML algorithm. Some of the categorical encoding techniques are described below.

**B.1. One-hot encoding**

This method converts an attribute with N possible categories into N distinct features. In the NSL-KDD dataset, the *protocol type* attribute has 3 possible values - Internet Message Control Protocol (ICMP), Transmission Control Protocol (TCP) and User Datagram Protocol (UDP). One-hot encoding converts this attribute into three feature columns as shown in Figure 11. It follows a 0/1 representation to indicate presence (indicated by 1) or absence (indicated by 0) of a value.

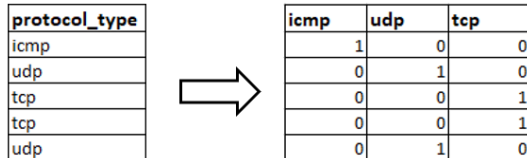


Figure 11. One-hot encoding used in Protocol Type attribute.

One-hot encoding has the advantage that it can convert ordinal and categorical features into orthogonal vector space [56]. It is one of the most used feature transformation techniques in machine learning experiments. The disadvantage of one-hot encoding appears when features with high cardinality are transformed. It might have the curse of dimensionality, also leading to a very sparse matrix. Also, there is a possibility to introduce bias in case the feature is distributed in an orderly fashion [57].

**B.2. Dummy coding**

Like one-hot encoding, dummy coding converts an attribute with N values to a feature set of N-1 values. The converted set of binary variables are called dummy variables. Figure 12 illustrates one-hot encoding and dummy coding applied to the same set of categorical records.

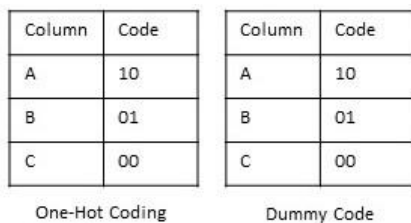


Figure 12. One-hot and dummy encoding used in the same dataset [12].

An advantage of dummy encoding is that it decreases the complexity and entropy complications. Dummy variable is of slightly better space complexity since it creates (k-1) dummy variables for k original variables. Whereas one hot encoding creates exactly k number of converted variables. The dummy variable trap is a camp term in machine learning preprocessing that refers to more than one independent categorical feature being multi-collinear [58][59].

**B.3. Effect coding**

While one-hot encoding and dummy coding use only 0 and 1 to encode categorical variables, effect coding sets values that sum to zero in the new feature set. As a result, negative values may also be generated in the encoded feature set. Effect coding is a preferred choice when there is an interaction of categorical variables in a dataset as it can provide reasonable estimates of main effect and of the interaction [13]. Effect coding has an advantage over dummy coding when there is an interaction between categorical variables [60]. The benefit is that the feature will achieve a reasonable estimate of both main effect and interaction using effect coding, which is efficient when training data is unbalanced [61].

**B.4. Hash encoding**

Hash encoding is appropriate for categorical variables that have many possible values. The method uses a hash function to map categorical values into numbers. Commonly used hash functions include Message Digest functions MD2, MD4, MD5 and Secure Hash Algorithms SHA0, SHA1, SHA2 and SHA3. The MD5 hash function is used by default [12]. Hash encoding returns a variable map with smaller dimensions than other encoding schemes, such as one-hot encoding or dummy coding. Figure 13 below illustrates the hash-encoding process:

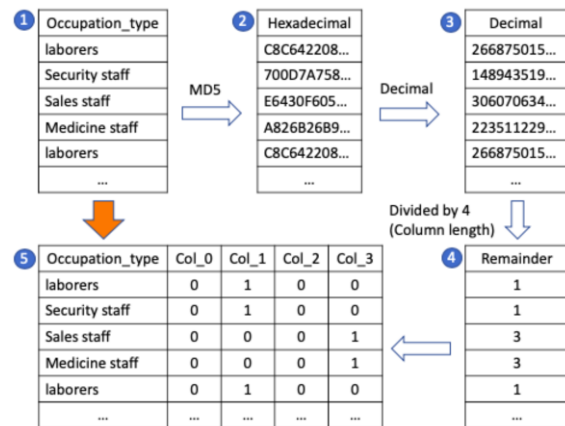


Figure 13. Hash Encoding Process [14].

Hash encoding is a fast and memory-efficient technique for categorical feature encoding. It has limitations such as inducing collisions and losing ordering information. Unlike the one-hot and dummy encoding, the method can handle unseen categorical values during interference by mapping them to a random numerical value [62]. Also, since it does not need a separate encoding dictionary, it is quite easy to deploy. Hash encoding can result in collisions, where different categorical values are mapped to the same numerical value. This can reduce the accuracy of the machine learning model and make it difficult to interpret the

importance of the feature. Hash encoding is considered unsuitable for categorical features with a small number of unique values, since the collision rate can be high, and the encoding may not be informative [63].

### B.5. BaseN encoding

BaseN encoding converts categorical variables into a consistently encoded feature set using a selected base, such as base 2 for binary encoding. The base or radix is the available number that can be used in different combinations to represent all values in a numbering system. A BaseN encoder encodes categorical values into arrays of their base-n representation.

BaseN encoding is a powerful technique for encoding categorical variables as numerical variables. BaseN encoding preserves ordinal relationships among categories, which can improve the interpretability and accuracy of the machine learning model [64]. It can reduce the number of features in the machine learning model, which improves computational efficiency and reduces risk of overfitting. Limitations are potential for noise or sparsity, inability to capture complex relationships, and compatibility with specific machine learning models. The suitability of BaseN encoding depends on the specific dataset and machine learning task at hand [65].

### B.6. Target encoding

Target encoding (also known as mean encoding) replaces a variable with a mean of the target value for that variable. Figure 14 provides an illustration. When the values of a categorical variable are already of a high volume, target encoding provides an advantage over other methods as it does not add extra dimensions to the dataset.

	feature	feature_label	feature_mean	target
0	Moscow	1	0.4	0
1	Moscow	1	0.4	1
2	Moscow	1	0.4	1
3	Moscow	1	0.4	0
4	Moscow	1	0.4	0
5	Tver	2	0.8	1
6	Tver	2	0.8	1
7	Tver	2	0.8	1
8	Tver	2	0.8	0

Figure 14. Target Encoding [15].

Target encoding can reduce the number of features in the machine learning model, which improves computational efficiency and reduces the risk of overfitting [66]. Target encoding is a powerful technique for encoding categorical variables as numerical variables based on the target variable. Limitations include potential for bias or overfitting, limited generalization to new data, and lack of compatibility with certain machine learning models. The suitability of Target

encoding depends on the specific dataset and machine learning task. It is recommended to use cross-validation to avoid overfitting when using target encoding [67].

### B.7. Label or ordinal encoding scheme

Ordinal categorical variables require that the order of the variable be preserved. For example, a road surface when a collision occurs might be categorized as dry, somewhat wet, or very wet so that the 3 values have an order that might provide additional insights. Ordinal encoding scheme aims at preserving this order when mapping values to numeric form. The method simply assigns each label a number (for example dry=1, somewhat wet=2 and very wet=3).

Label encoding preserves ordinal relationships among categories, which improves the interpretability and accuracy of the model [68]. It can reduce the number of features in the machine learning model, improving computational efficiency, and reducing risk of overfitting. Limitations are potential for bias or arbitrary ranking, inability to capture complex relationships, and lack of compatibility with specific machine learning models. The suitability of Label encoding depends on the specific dataset and machine learning task addressed. It is recommended to use label encoding only when the categorical variable has a meaningful order, or the number of categories is small [69].

### C. Feature Selection

The complexity of ML models increases with the dimensionality of the dataset. Predictive models often fail to achieve high accuracy because of inadequate analyses directed to feature selection. Selecting the most important and significant features reduces the complexity of the model and can also increase the prediction performance [37][38]. Multiple approaches are available and effective for reducing the feature set. Prominent ones are described below.

**Filter Methods.** In our work, we choose feature selection methods that apply to situations with a categorical output, such as whether an input vector is an attack or not. The filter methods eliminate features independently of the ML method used. A univariate feature filter evaluates the importance of single features using univariate statistical tests. Each feature is paired with the target to evaluate statistical significance between them. The analysis of variance or ANOVA F-test is widely used. The F-test calculates the ratio between variance values [31]. The resultant measures of the relative importance of individual features provide a tool for determining features that are unnecessary or of little importance.

The filter method is effective for scalability, especially for high-dimensional datasets with a large number of features [70]. It improves accuracy and interpretability by removing irrelevant and redundant features. It also helps to reduce overfitting and improve generalization. The disadvantages of the filter method are the lack of modeling interaction and/or non-linear relationships between features and the target

variables [79]. The method limits the accuracy and predictive power to a certain threshold [71].

**Wrapper Methods.** Wrapper methods directly evaluate combinations of features by running the ML model restricted to the set of candidate features. Taken to an extreme, all combinations would be evaluated, an impossible task in practice. Thus, various search space approaches are employed. Forward search iteratively adds promising feature vectors one by one to build a feature set. Backward search starts with all features and successively eliminates those that perform poorly. Many approaches based on optimization techniques are available [32][33]. The self-driving car work basically follows a forward selection approach based on both advance knowledge about the importance of certain features from analytics on the data itself along with test runs of the ML model. Heavy computational load and possibilities of overfitting are potential drawbacks [34].

The Wrapper method is efficient in capturing complex relationships and interactions between features and targets. It can handle non-linear relationships which is often impossible for other feature selection techniques [72]. Like the filter method, it improves the interpretability and generalization ability of the model. Wrapper methods are computationally expensive, especially for high dimensional data. Also, the wrapper method can overfit the model to the training data, especially if the search space of feature subsets is large or the evaluation metric is not carefully chosen [73].

**Embedded Methods.** Embedded methods utilize mathematical information that is available during the training of the model to determine the relative importance of features. In some sense embedded methods mitigate the drawbacks of the filter and wrapper methods but retain their strengths. When implemented carefully, they are not prone to overfitting [35]. The XGBoost technique produces an importance score for each attribute that is used to identify those that can be confidently eliminated [36]. In applications like intrusion detection, having many attributes presents a huge computational burden. The embedded methods are highly successful in greatly reducing the features needed in intrusion detection ML work.

The Embedded method handles non-linear relationships and interactions among features and the target variable, a capability that may not be possessed by the filter method or many other feature selection techniques [74]. The method is a powerful technique for feature selection that can capture complex relationships and interactions among features and the target variable. Limitations include high computational cost, lack of compatibility with certain models, and potential for overfitting or underfitting. The suitability of embedded method depends on the specific dataset, machine learning model, and optimization algorithm [75].

## VI. CONCLUSION AND FUTURE WORK

Machine learning is now a well-established and effective approach in many domains. When using machine learning approaches in practice, issues arise concerning choices for which type of model to use, parameters to choose and tune, data management, feature engineering and selection. We address many of these issues in the context of applications to self-driving cars and intrusion detection. The applications are of high importance in inter-related areas of cybersecurity, trust, risk, safety, reliability, autonomy, and anti-autonomy. For the studies concerning self-driving cars, we present uses for reinforcement learning and supervised learning that reveal circumstances under which the autonomous vehicle makes decisions to take actions that result in collisions. For intrusion detection, model choices and computational results are presented that result in excellent values for performance metrics. Included are reports of recall metrics that indicate that it is possible to use machine learning methods that flag nearly all potentially dangerous attack vectors. The data-centric and feature engineering challenges are extensive and detailed, but addressable. We describe approaches to addressing these challenges. Results reveal several implications for needs for next steps in research. New frontiers include methods that can be deployed in real-time, automate feature engineering, choose, and extract features dynamically, and simultaneously support optimization of multiple performance evaluation metrics.

## REFERENCES

- [1] K. E. Nygard, M. Ahsan, A. Rastogi, and R. Satyal, "Data and Feature Engineering Challenges in Machine Learning," in *ADVCOMP 2022, The Sixteenth International Conference on Advanced Engineering Computing and Applications in Sciences*, pp. 1–10, November 2022.
- [2] T. Hey, S. Tansley, and K. Tolle, "The Fourth Paradigm: Data-Intensive Scientific Discovery," Redmond, Washington: Microsoft Research, 2009.
- [3] A. Rastogi, "Trust and Anti-Autonomy Modelling of Autonomous Systems," Order No. 28255688, North Dakota State University, Ann Arbor, 2020.
- [4] R. Meiri and J. Zahavi, "Using simulated annealing to optimize the feature selection problem in marketing applications," in *European Journal of Operational Research*, vol. 171, no. 3, pp. 842–858, 2006.
- [5] M. Lombardi and M. Milano, "Boosting combinatorial problem modeling with machine learning," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pp. 5472–5478, 2018.
- [6] M. Charikar, V. Guruswami, R. Kumar, S. Rajagopalan, and A. Sahai, "Combinatorial feature selection problems," in *Proceedings 41st Annual Symposium on Foundations of Computer Science*, pp. 631–640, 2000.
- [7] KDD Cup 1999 Data, [kdd.ics.uci.edu](http://kdd.ics.uci.edu) [Online]. Available from: <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html> [Accessed: 2022.08.22].

- [8] A Deeper Dive into the NSL-KDD Data Set. Medium [Online]. Available from: <https://towardsdatascience.com/a-deeper-dive-into-the-nsl-kdd-data-set-15c753364657> [Accessed: 2022.08.22].
- [9] L. Dhanabal and S.P. Shantharajah, "A Study on NSL-KDD Dataset for Intrusion Detection System Based on Classification Algorithms," In *International Journal of Advanced Research in Computer and Communication Engineering*, vol. 4, no. 6, pp. 446-452, June 2015.
- [10] S. Choudharya and N. Kesswani, "Analysis of KDD-Cup'99, NSL-KDD and UNSW-NB15 Datasets using Deep Learning in IoT," In *Procedia Computer Science*, vol. 167, pp. 1561-1573, 2020.
- [11] A. Divekar, M. Parekh, V. Savla, R. Mishra, and M. Shirole, "Benchmarking datasets for Anomaly-based Network Intrusion Detection: KDD CUP 99 alternatives," In *2018 IEEE 3rd International Conference on Computing, Communication and Security (ICCCS)*, IEEE, pp. 1-8, 2018.
- [12] 8 Categorical Data Encoding Techniques to Boost your Model in Python!, Analytics Vidhya [Online]. Available from: <https://www.analyticsvidhya.com/blog/2020/08/types-of-categorical-data-encoding/> [Accessed: 2022.08.22].
- [13] Introduction to SAS, UCLA: Statistical Consulting Group. [Online]. Available from: <https://stats.idre.ucla.edu/sas/modules/sas-learning-moduleintroduction-to-the-features-of-sas/> [Accessed: 2022.08.22]
- [14] A Data Scientist's Toolkit to Encode Categorical Variables to Numeric. [Online]. Available from: <https://towardsdatascience.com/a-data-scientists-toolkit-to-encode-categorical-variables-to-numeric-d17ad9fae03f> [Accessed: 2022.08.22]
- [15] Improve your classification models using Mean /Target Encoding. [Online]. Available from: <https://medium.com/datadriveninvestor/improve-your-classification-models-using-mean-target-encoding-a3d573df31e8> [Accessed: 2022.08.22]
- [16] N. Moustafa and J. Slay, "UNSW-NB15: a comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set)," In *2015 military communications and information systems conference (MilCIS)*, IEEE, pp. 1-6, 2015.
- [17] N. Japkowicz and S. Stephen, "The class imbalance problem: A systematic study," In *Intelligent data analysis*, vol. 6, no. 5, pp. 429-449, 2002.
- [18] R. Gomes, M. Ahsan, and A. Denton, "Random Forest Classifier in SDN Framework for User-Based Indoor Localization," In *2018 IEEE International Conference on Electro/Information Technology (EIT)*, IEEE, pp. 0537-0542, 2018.
- [19] A. Ali, S. M. Shamsuddin, and A. L. Ralescu, "Classification with class imbalance problem: A review," In *Int. J. Adv. Soft Comput. its Appl.*, vol. 5, no. 3, 2013.
- [20] A. Ghazikhani, H. S. Yazdi, and R. Monsefi, "Class imbalance handling using wrapper-based random oversampling," In *20th Iranian Conference on Electrical Engineering (ICEE2012)*, IEEE, pp. 611-616, 2012.
- [21] Z. Zheng, Y. Cai, and Y. Li, "Oversampling method for imbalanced classification," In *Computing and Informatics*, vol. 34, no. 5, pp. 1017-1037, 2015.
- [22] A. M. Denton, M. Ahsan, D. Franzen, and J. Nowatzki, "Multi-scalar analysis of geospatial agricultural data for sustainability," In *2016 IEEE International Conference on Big Data (Big Data)*, IEEE, pp. 2139-2146, 2016.
- [23] N. V. Chawla, N. Japkowicz, and A. Kotcz, "Special Issue on Learning from Imbalanced Data Sets," In *ACM SIGKDD explorations newsletter*, vol. 6, no. 1, pp. 1-6, 2004.
- [24] M. Ahsan, R. Gomes, and A. Denton, "SMOTE Implementation on Phishing Data to Enhance Cybersecurity," In *2018 IEEE International Conference on Electro/Information Technology (EIT)*, IEEE, pp. 0531-0536, 2018.
- [25] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic minority over-sampling technique," In *Journal of artificial intelligence research*, vol. 16, pp. 321-357, 2002.
- [26] H. Han, W. Y. Wang, and B. H. Mao, "Borderline-SMOTE: A new over-sampling method in imbalanced data sets learning," In *International conference on intelligent computing*, pp. 878-887. Springer, Berlin, Heidelberg, 2005.
- [27] H. Hairani, K. E. Saputro, and S. Fadli, "K-means-SMOTE untuk menangani ketidakseimbangan kelas dalam klasifikasi penyakit diabetes dengan C4.5, SVM, dan naive Bayes," In *J. Teknol. dan Sist. Komput.*, vol. 8, no. 2, pp. 89-93, 2020.
- English – H. Hairani, K. E. Saputro, and S. Fadli, "K-means-SMOTE to trate class imbalance in the classification of diabetes with C4.5, SVM, and Naive Bayes," In *J. Teknol. dan Sist. Komput.*, vol. 8, no. 2, pp. 89-93, 2020.
- [28] Y. Tang, Y. Q. Zhang, and N. V. Chawla, "SVMs modeling for highly imbalanced classification," In *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 39, no. 1, pp. 281-288, 2008.
- [29] H. He, Y. Bai, E. A. Garcia, and S. Li, "ADASYN: Adaptive synthetic sampling approach for imbalanced learning," In *2008 IEEE international joint conference on neural networks (IEEE world congress on computational intelligence)*, IEEE, pp. 1322-1328, 2008.
- [30] O. Özdemir, M. Batar, and A. H. Işık , "Churn Analysis with Machine Learning Classification Algorithms in Python," In *The International Conference on Artificial Intelligence and Applied Mathematics in Engineering*, pp. 844-852. Springer, Cham, 2019.
- [31] N. O. F. Elssied, O. Ibrahim, and A. H. Osman, "A novel feature selection based on one-way ANOVA F-test for e-mail spam classification," In *Research Journal of Applied Sciences, Engineering and Technology*, vol. 7, no. 3, pp. 625-638, 2014.
- [32] M. Ahsan, R. Gomes, and A. Denton, "Application of a convolutional neural network using transfer learning for tuberculosis detection," In *2019 IEEE International Conference on Electro Information Technology (EIT)*, IEEE, pp. 427-433, 2019.
- [33] A. Mustaqeem, S. M. Anwar, M. Majid, and A. R. Khan, "Wrapper method for feature selection to classify cardiac arrhythmia," In *2017 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, IEEE, pp. 3656-3659, 2017.
- [34] M. Ahsan and K. Nygard, "Convolutional Neural Networks with LSTM for Intrusion Detection," In *The 35th International*

- Conference on Computers and Their Applications (CATA 2018)*, vol. 69, pp. 69-79, 2020.
- [35] H. Liu, M. Zhou, and Q. Liu, "An embedded feature selection method for imbalanced data classification," In *IEEE/CAA Journal of Automatica Sinica*, vol. 6, no. 3, pp. 703-715, 2019.
- [36] Y. Wang and X. S. Ni, "A xgboost risk model via feature selection and bayesian hyper-parameter optimization," In *arXiv preprint arXiv:1901.08433*, 2019.
- [37] M. Chowdhury, J. Tang, and K. Nygard, "An artificial immune system heuristic in a smart grid," In *The 28th International Conference on Computers and Their Applications*, 2013.
- [38] M. Chowdhury and K. Nygard, "Machine learning within a con resistant trust model," In *The 33rd International Conference on Computers and Their Applications (CATA 2018)*, pp. 9-14, 2018.
- [39] A. Rastogi and K. Nygard, "Threat and alert analytics in autonomous vehicles," In *The 35th International Conference on Computers and Their Applications (CATA 2018)*, vol. 69, pp. 48-59, 2020.
- [40] Home - Keras Documentation, Keras.io, 2020. [Online]. Available: <https://keras.io/>. [Accessed: 2022.08.16].
- [41] "TensorFlow," TensorFlow, 2020. [Online]. Available: <https://www.tensorflow.org/>. [Accessed: 2022.08.16].
- [42] "tensorflow/tensorflow," GitHub, 2020. [Online]. Available: <https://github.com/tensorflow/tensorflow>. [Accessed: 2022.08.16].
- [43] M. Ahsan, N. Rifat, M. Chowdhury, and R. Gomes, "Detecting Cyber Attacks: A Reinforcement Learning Based Intrusion Detection System," In *2022 IEEE International Conference on Electro Information Technology (eIT)*, IEEE, pp. 461-466, 2022.
- [44] M. Ahsan, K. E. Nygard, R. Gomes, M. M. Chowdhury, N. Rifat, and J. F. Connolly, "Cybersecurity Threats and Their Mitigation Approaches Using Machine Learning—A Review," In *Journal of Cybersecurity and Privacy*, vol. 2, no. 3, pp. 527-555, 2022.
- [45] N. I. Rifat, "Feature Engineering on the Cybersecurity Dataset for Deployment on Software Defined Network," 2020.
- [46] Y. Berra, "A Quote By Yogi Berra," 2022. [Online]. Available: <https://www.goodreads.com/quotes/261863-it-s-tough-to-make-predictions-especially-about-the-future>. [Accessed: 2022.08.22].
- [47] M. Ahsan, N. Rifat, M. Chowdhury, and R. Gomes, "Intrusion Detection for IoT Network Security with Deep Neural Network," In *2022 IEEE International Conference on Electro Information Technology (eIT)*, IEEE, pp. 467-472, 2022.
- [48] P. Sajda, A. Gerson, K. R. Muller, B. Blankertz, and L. Parra, "A data analysis competition to evaluate machine learning algorithms for use in brain-computer interfaces," In *IEEE Transactions on neural systems and rehabilitation engineering*, 11(2), pp. 184-185, 2003.
- [49] A. Fleury-Charles, M. M. Chowdhury, and N. Rifat, "Data Breaches: Vulnerable Privacy," In *2022 IEEE International Conference on Electro Information Technology (eIT)*, Mankato, MN, USA, pp. 538-543, 2022.
- [50] Z. Li, X. Ma, and H. Xin, "Feature engineering of machine-learning chemisorption models for catalyst design," *Catalysis today*, vol. 280, pp. 232-238, 2017.
- [51] D. K. Davis, M. M. Chowdhury, and N. Rifat, "Password Security: What Are We Doing Wrong?," In *2022 IEEE International Conference on Electro Information Technology (eIT)*, Mankato, MN, USA, 2022, pp. 562-567.
- [52] R. Vanness, M. M. Chowdhury, and N. Rifat, "Malware: A Software for Cybercrime," In *2022 IEEE International Conference on Electro Information Technology (eIT)*, Mankato, MN, USA, pp. 513-518, 2022.
- [53] Y. Xu, K. Hong, J. Tsujii, and E.I.C. Chang, "Feature engineering combined with machine learning and rule-based methods for structured information extraction from narrative clinical discharge summaries," In *Journal of the American Medical Informatics Association*, vol. 19, no. 5, pp. 824-832, 2012.
- [54] M.Ahsan, R. Gomes, M.M. Chowdhury, and K. E. Nygard, "Enhancing machine learning prediction in cybersecurity using dynamic feature selector," In *Journal of Cybersecurity and Privacy*, vol. 1, no. 1, pp. 199-218, 2021.
- [55] M.K. Ahsan, "Increasing the Predictive Potential of Machine Learning Models for Enhancing Cybersecurity," Order No. 28416096, North Dakota State University, Ann Arbor, 2021.
- [56] M. K. Dahouda and I. Joe, "A deep-learned embedding technique for categorical features encoding," In *IEEE Access*, vol. 9, pp. 114381-114391, 2021.
- [57] B. Gu and Y. Sung, "Enhanced reinforcement learning method combining one-hot encoding-based vectors for CNN-based alternative high-level decisions," In *Applied Sciences*, vol. 11, no.3, p. 1291, 2021.
- [58] I. Mukherjee, N.K. Sahu, and S.K. Sahana, "Simulation and modeling for anomaly detection in IoT network using machine learning," In *International Journal of Wireless Information Networks*, pp. 1-17, 2022.
- [59] J.M. Jerez, I. Molina, P.J. García-Laencina, E. Alba, N. Ribelles, M. Martín, and L. Franco, "Missing data imputation using statistical and machine learning methods in a real breast cancer problem," In *Artificial intelligence in medicine*, vol. 50, no.2, pp. 105-115, 2010.
- [60] R.S. Sutton, "Generalization in reinforcement learning: Successful examples using sparse coarse coding," In *Advances in neural information processing systems*, vol. 8, 1995.
- [61] C.H. Chien, C.C. Chang, S.H. Lin, C.W. Chen, Z.H. Chang, and Y.W. Chu, "N-GlycoGo: predicting protein N-glycosylation sites on imbalanced data sets by using heterogeneous and comprehensive strategy," In *IEEE Access*, vol. 8, pp. 165944-165950, 2020.
- [62] Z. Cao, M. Long, J. Wang, and P.S. Yu, "Hashnet: Deep learning to hash by continuation," In *Proceedings of the IEEE international conference on computer vision*, pp. 5608-5617, 2017.
- [63] I. Lopez-Arevalo, E. Aldana-Bobadilla, A. Molina-Villegas, H. Galeana-Zapién, V. Muñoz-Sanchez, and S. Gausin-Valle, "A memory-efficient encoding method for processing mixed-type data on machine learning," In *Entropy*, vol. 22, no.12, p. 1391, 2020.
- [64] J.M Springer, C.S. Strauss, A.M. Thresher, E. Kim, and G.T. Kenyon, "Classifiers based on deep sparse coding architectures

- are robust to deep learning transferable examples,” In *arXiv preprint arXiv:1811.07211*, 2018.
- [65] S. Naseer, Y. Saleem, S. Khalid, M.K. Bashir, J. Han, M.M. Iqbal, and K. Han, “Enhanced network anomaly detection based on deep neural networks,” In *IEEE Access*, vol. 6, pp. 48231-48246, 2018.
- [66] N. Nazyrova, T.J. Chausalet, and S.Chahed, “Machine Learning models for predicting 30-day readmission of elderly patients using custom target encoding approach,” In *Computational Science–ICCS 2022: 22nd International Conference Proceedings*, vol. III, pp. 122-136, June 2022.
- [67] D. Silhavy, C. Krauss, A. Chen, A.T. Nguyen, C. Müller, S. Arbanowski, S. Steglich, and L. Bassbouss, “Machine learning for per-title encoding,” In *SMPTE Motion Imaging Journal*, vol. 131, no.3, pp. 42-50, 2022.
- [68] A. Mottini and R. Acuna-Agost, “Relative label encoding for the prediction of airline passenger nationality,” In *2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW)*, pp. 671-676, 2016.
- [69] C. Song, T. Ristenpart, and V. Shmatikov, “Machine learning models that remember too much,” In *Proceedings of the 2017 ACM SIGSAC Conference on computer and communications security*, pp. 587-601, 2017.
- [70] S. Khalid, T. Khalil, and S. Nasreen, “A survey of feature selection and feature extraction techniques in machine learning,” In *2014 science and information conference*, pp. 372-378, August 2014.
- [71] A. Zheng and A. Casari, “Feature engineering for machine learning: principles and techniques for data scientists,” In O'Reilly Media, Inc., 2018.
- [72] S.M. Kasongo and Y. Sun, “A deep learning method with wrapper based feature extraction for wireless intrusion detection system,” In *Computers & Security*, vol. 92, pp. 101752, 2020.
- [73] M. F. Rafique, M. Ali, A. S. Qureshi, A. Khan, and A. M. Mirza, “Malware classification using deep learning based feature extraction and wrapper based feature selection technique,” In *arXiv preprint arXiv:1910.10958*, 2019.
- [74] J. Zhou, L. Liu, W. Wei, and J. Fan, “Network representation learning: from preprocessing, feature extraction to node embedding,” In *ACM Computing Surveys (CSUR)*, vol. 55, no. 2, pp. 1-35, 2022.
- [75] S. Wang, J. Arroyo, J.T. Vogelstein, and C. E. Priebe, “Joint embedding of graphs,” In *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no.4, pp. 1324-1336, 2019.
- [76] Introduction to Various Reinforcement Learning Algorithms. Medium [Online]. Available from: <https://medium.com/towards-data-science/introduction-to-various-reinforcement-learning-algorithms-i-q-learning-sarsa-dqn-ddpg-72a5e0cb6287> [Accessed: 2023.03.05].
- [77] A. Rastogi and K. E. Nygard, “Trust Issues in Cybersecurity and Autonomy,” In *27<sup>th</sup> International Conference on Software Engineering and Data Engineering*, 2018.
- [78] U. S. Shanthamallu, A. Spanias, C. Tepedelenioglu, and M. Stanley, “A brief survey of machine learning methods and their sensor and IoT applications,” In *2017 8th International Conference on Information, Intelligence, Systems & Applications (IISA)*, pp. 1-8, 2017.
- [79] S.M. Kasongo and Y. Sun, “A deep learning method with filter based feature engineering for a wireless intrusion detection system,” In *IEEE Access*, vol. 7, pp. 38597-38607, 2019.