Editor-in-Chief   Professor Bekir Karlik

# INTERNATIONAL JOURNAL OF ARTIFICIAL INTELLIGENCE AND EXPERT SYSTEMS (IJAE)

# INTERNATIONAL JOURNAL OF ARTIFICIAL INTELLIGENCE AND EXPERT SYSTEMS (IJAE)

**VOLUME 10, ISSUE 2, 2021**

**EDITED BY
DR. NABEEL TAHIR**

# INTERNATIONAL JOURNAL OF ARTIFICIAL INTELLIGENCE AND EXPERT SYSTEMS (IJAE)

**CSC Publishers, 2021**

# EDITORIAL BOARD

# TABLE OF CONTENTS

## Pages

# EDITORIAL PREFACE

The International Journal of Artificial Intelligence and Expert Systems (IJAE) is an effective medium for interchange of high quality theoretical and applied research in Artificial Intelligence and Expert Systems domain from theoretical research to application development. This is the *Second* Issue of Volume *Ten* of IJAE. The Journal is published bi-monthly, with papers being peer reviewed to high international standards. IJAE emphasizes on efficient and effective Artificial Intelligence, and provides a central for a deeper understanding in the discipline by encouraging the quantitative comparison and performance evaluation of the emerging components of Expert Systems. IJAE comprehensively cover the system, processing and application aspects of Artificial Intelligence. Some of the important topics are AI for Service Engineering and Automated Reasoning, Evolutionary and Swarm Algorithms and Expert System Development Stages, Fuzzy Sets and logic and Knowledge-Based Systems, Problem solving Methods Self-Healing and Autonomous Systems etc.

The initial efforts helped to shape the editorial policy and to sharpen the focus of the journal. Started with Volume 10, 2021, IJAE will be appearing with more focused issues related to artificial intelligence and expert system research. Besides normal publications, IJAE intend to organized special issues on more focused topics. Each special issue will have a designated editor (editors) – either member of the editorial board or another recognized specialist in the respective field.

IJAE give an opportunity to scientists, researchers, and vendors from different disciplines of Artificial Intelligence to share the ideas, identify problems, investigate relevant issues, share common interests, explore new approaches, and initiate possible collaborative research and system development. This journal is helpful for the researchers and R&D engineers, scientists all those persons who are involve in Artificial Intelligence and Expert Systems in any shape.

Highly professional scholars give their efforts, valuable time, expertise and motivation to IJAE as Editorial board members. All submissions are evaluated by the International Editorial Board. The International Editorial Board ensures that significant developments in image processing from around the world are reflected in the IJAE publications.

IJAE editors understand that how much it is important for authors and researchers to have their work published with a minimum delay after submission of their papers. They also strongly believe that the direct communication between the editors and authors are important for the welfare, quality and wellbeing of the Journal and its readers. Therefore, all activities from paper submission to paper publication are controlled through electronic systems that include electronic submission, editorial panel and review system that ensures rapid decision with least delays in the publication processes.

To build its international reputation, we are disseminating the publication information through Google Books, Google Scholar, Directory of Open Access Journals (DOAJ), Open J Gate, ScientificCommons, Docstoc and many more. Our International Editors are working on establishing ISI listing and a good impact factor for IJAE. We would like to remind you that the success of our journal depends directly on the number of quality articles submitted for review. Accordingly, we would like to request your participation by submitting quality manuscripts for review and encouraging your colleagues to submit quality manuscripts for review. One of the great benefits we can provide to our prospective authors is the mentoring nature of our review process. IJAE provides authors with high quality, helpful reviews that are shaped to assist authors in improving their manuscripts.

**Editorial Board Members**
International Journal of Artificial Intelligence and Expert Systems (IJAE)

Carina Dantas, Karolina Mackiewicz, Valentina Tageo, Giulio Jacucci, Diana Guardado, Sofia Ortet, Iraklis Varlamis, Michail Maniadakis, Eva de Lera, João Quintas, Otilia Kocsis & Charalampos Vassiliou

# Benefits and Hurdles of AI In The Workplace – What Comes Next?

**Carina Dantas**                                            *carina@echalliance.com*
*European Connected Health Alliance*
*Coimbra, 3030 Portugal*


**Karolina Mackiewicz**                                   *karolina@echalliance.com*
*European Connected Health Alliance*
*Finland*


**Valentina Tageo**                                        *valentina@echalliance.com*
*European Connected Health Alliance*
*Barcelona, Postal code, Spain*


**Giulio Jacucci**                                         *giulio.jacucci@helsinki.fi*
*Department of Computer Science*
*University of Helsinki, Finland*


**Diana Guardado**                               *dianaguardado@caritascoimbra.pt*
*Cáritas Diocesana de Coimbra*
*Coimbra, 3030 Portugal*


**Sofia Ortet**                                       *sofiaortet@caritascoimbra.pt*
*Cáritas Diocesana de Coimbra*
*Coimbra, 3030 Portugal*


**Iraklis Varlamis**                                            *varlamis@hua.gr*
*Department of Informatics and Telematics*
*Harokopio Univercity of Athens, Greece*
*& ICS FORTH, Athens, 177 78, Greece*


**Michail Maniadakis**                                    *mmaniada@ics.forth.gr*
*ICS, Foundation for Research and Technology*
*Heraklion, 700 13, Greece*


**Eva de Lera**                                        *Eva@raisingthefloor.org*
*Raising the Floor - International Association*
*Geneva, CH-1218, Switzerland*


**João Quintas**                                               *jquintas@ipn.pt*
*Automatics Laboratory*
*Instituto Pedro Nunes*
*Coimbra, 3030-199, Portugal*


**Otilia Kocsis**                                              *okocsis@bok.gr*
*Department of Electrical and Computer Engineering*
*University of Patras*
*Patras, 26500, Greece*


**Charalampos Vassiliou**                                    *cvassiliou@byte.gr*
*Project Coordinator*
*BYTE*
*117 41 Athens, Greece*

Carina Dantas, Karolina Mackiewicz, Valentina Tageo, Giulio Jacucci, Diana Guardado, Sofia Ortet, Iraklis Varlamis, Michail Maniadakis, Eva de Lera, João Quintas, Otilia Kocsis & Charalampos Vassiliou

**Abstract**

During the last few years, there has already been a solid discussion and political will, transversal to most European countries, on the need to invest in prevention, promoting healthier living environments and person-centred digital solutions. In short, it seems that consensus on the need to move forward and invest in wellbeing and quality of life was achieved.

During the COVID-19 pandemic and the confinement measures it implied, many services had to be closed; teleworking suddenly became the rule and many families stayed at their homes, with the children in remote classes, some without sufficient equipment or the most adequate digital tools available. Digital services, if implemented correctly, can be the right tools to address many of these challenges. The importance of implementing them correctly increases in the current context of accelerated Digital Transformation, where many are looking towards Artificial Intelligence (AI) as a means to help people to cope with the increasing number of digitized work. We are assisting to a gear-shift in the current digital revolution, as now we better understand how it could have been helpful, if already embedded in daily life.

COVID-19 generated severe consequences for the working context, with effects on physical and mental health and wellbeing, and with trends such as teleworking coming to stay. Organizations and individuals working on AI can play a great role in providing solutions, not only during this emergency period, but also in the long-term perspective, and not only for office workers but in more traditional industries as well. Thus, the COVID-19 pandemic is a driver for the digital revolution in the workplace across many levels. However, inequalities persist and their impact on universal access to the digital world is enormous. Moreover, several other challenges come from the use of artificial intelligence in the workplace.

This paper addresses how technology applied to the work environment can be leveraged to respond to the emerging challenges raised by COVID-19. It also provides reflections on the main opportunities and challenges that the use of AI solutions in the workplace imply, suggesting measures or recommendations to tackle them, towards a concerted approach to AI, integrating the policy agenda with the implementation strategy.

**Keywords:** AI, Living and Working Environments, Digital Transformation, COVID19, Inclusion.

---

## 1. INTRODUCTION

Digital transformation is a latest trend, related to the transformation in business and operations by utilizing digital technologies. It is considered a major development ankle for corporations, supporting them against competition and enhancing their mid-term viability [1]. Usually, the deployment of digital technologies for supporting different aspects of a business (from sales [2] and marketing, to everyday operations and financial management), is based on the vision and decision making of the managers, who are responsible for defining and monitoring the organisation's long-term strategy. Furthermore, the application of new technological trends (e.g. teleworking) in working environments is a gradual and time-consuming process [3], especially in the most traditional ones (e.g. banking, public sector), requiring a long transition period, during which the employees need to be trained to acquire the sufficient digital literacy, potentiating the gradual acceptance of the changes.

However, the recent COVID-19 pandemic and its results proved to be a greater disruptor for the existing working environments [4]. The enforced lockdowns and physical distancing policies made it almost impossible for companies to maintain their full force at their premises, and had to proceed fast to an organizational transformation, from office settings to a fully remote format, in order to keep their personnel safe and healthy and their business running. Despite the fact that many employees were also allowed previously to work from home for a minimum amount of their time (e.g. four times per month), such policies have never applied before at such a great extent.

Carina Dantas, Karolina Mackiewicz, Valentina Tageo, Giulio Jacucci, Diana Guardado, Sofia Ortet, Iraklis Varlamis, Michail Maniadakis, Eva de Lera, João Quintas, Otilia Kocsis & Charalampos Vassiliou

In order to prevent employers from interrupting their business activities, as employees from losing their income, COVID19 pandemic has been triggering remote work in a way that, probably, may never be completely reversed to its original state [5]. Modern computer technology and the digital transformation is actually the main responsible for allowing the majority of office workers to continue with their work, even if not in their usual workplace.

It is also important to note that even though digital solutions were initially welcomed as the tools that would facilitate processes and remove demands, pressure and risks, in the workplace amongst many other environments, the resulting effect has been quite the opposite, it has increased the demand of new skills (e.g digital literacy), it has increased the expectations of outcomes from employees, and also created a 24/7 environment (mostly through mobile technologies) which, in many cases, demanded a transition from a scheduled limited availability, to a constant availability [6]. Such increase in expectations, availability and technological requirements has put a strain in the quality of life of the workforce, in particular in the case of older workers [7]. It is then critical to learn about what is happening during this digital-everywhere-all the time transitions, and generate the necessary insights to guarantee and help improve overall quality of work, and of life.

In the present article, the authors departed from several European projects reports, implementation experiments, events and research discussion, as well as an extensive literature search undertaken under these different initiatives, to advance with some considerations regarding the topic of AI in the workplace and the effects of COVID19 in the challenges to overcome, mainly based on inductive reasoning, supported by the shared opinions of a wide and multidisciplinary team of experts.

## 2. USE OF DIGITAL TOOLS IN THE WORKPLACE

Technology proved to be the greatest enabler for successfully completing this transformational phase of the workplace during the COVID-19 pandemic [8]. By leveraging the capabilities of services already applied in the office environments, and with the addition of a limited number of novel ones, a virtual office space was formulated where all employees may interact, collaborate and participate in company's everyday activities while being physically located at their homes [9].

It is important to note that not all organizations and people could adapt to these new needs and environments, in some cases due to lack of budget and resources (i.e. hardware, Wi-Fi), skills or adaptive systems (i.e. the inexistence of assistive technologies) for this new hybrid-online working.

When integrating AI systems in the workplace, all of the above is to be considered, to ensure a realistic input of information into these systems, to ensure that its benefits are greater, and not limited.

### 2.1 Benefits and Hurdles

For instance, tele-conference services, like Teams, Skype, Zoom, GoToMeeting, Webex, along with more informal channels (e.g. viber, WhatsApp, Signal, Slack, etc.) tackled the communication needs. Cloud-based repositories like Dropbox, Google Drive and SharePoint enhanced the exchange of files, documents, and information. Online office tools, like Google Suite or Office 365 allowed for the simultaneous access and editing of their content. VPN tools enabled the access to internal information stored in a company's server. Last but not least, collaboration and project management tools like Trello, Redmine or JIRA made possible the monitoring of a project progress and the efficient task allocation. The list of digital tools is long and covers all operations of a company or other type of organisation, from human resources and high-level management to IT services and marketing and sales, ensuring the efficient collaboration in a fully remote environment.

However, the adoption of technology cannot solely make the difference leading to the aforementioned transition, as several non-technological challenges need to be also tackled. For

Carina Dantas, Karolina Mackiewicz, Valentina Tageo, Giulio Jacucci, Diana Guardado, Sofia Ortet, Iraklis Varlamis, Michail Maniadakis, Eva de Lera, João Quintas, Otilia Kocsis & Charalampos Vassiliou

instance, the vast majority of companies from smaller ones to larger multinational organisations have established procedures and processes that formalized the operations and set the framework for running their businesses. Therefore, in order to make it fully running in a remote environment, the procedures should be adjusted accordingly. Modifying or adjusting these, in order to exploit the available digital tools may prove harder than learning to use the technology itself. Indeed, the modification in aspects that were successful and became part of an organisation's comfort zone (business-as-usual) may create a high resistance to change, limiting the impact of the digital intervention. Issues such as data governance, ethics and privacy assume a much more important role when the work increases its digitalisation.

Similar to the procedural changes, another challenge that needs to be tackled, is the change in a company's nature. The transition in a remote setting may have been proven easier for companies already in the technology domain (e.g. software houses) or companies using a great deal of digital tools. But what would happen in a more "traditional" business (e.g. food production)? In this case, the adoption of physical distancing policies within the companies' promises, the enforcement of strict hygiene measures and the implementation of shift-work programs, will have a better effect for such organizations, whose premises are core elements of their business. Even here, technology may play a supporting but substantial role in the effective employee management and the monitoring of physical/social distancing policies.

Last but not least, one of the greatest challenges during the COVID-19 period is how to keep the employees healthy and safe. As the employees are the heart of a company, making sure that they will not get sick is crucial for keeping the business running. In addition, this period was proven to be extremely stressful, the consequences of which in personal health have been measured yet. Therefore, maintaining the good mental state of the employees is of equal importance. Preserving the good physical and mental health of the employees does not only concern the older ones, who may belong to a vulnerable group and consequently are more susceptible to COVID-19, but also affect younger ones who may be living with their parents or families, are more concerned about the future and may experience higher stressful conditions. Technology interventions in this case, may help to monitor their health and to ensure the good level of their mental and physical wellness.

Technology can have a role into workers environment also when they return to the workplace when restrictions are relaxed. In this context tracing technologies and behavior change approaches can support workers in developing working practices that are safe and minimize the chance of contracting the virus. It has come to light that maintaining a healthy environment is as much a social acceptance and behavior change problem as much as it is a problem of protection equipment and vaccines.

There seems to be some sort of lack of control of the use of technology, as it has entered people's personal and professional lives, increasing the number of online interactions in a given day, with colleagues, vendors, clients, friends and family members. It appears there is a need to spread awareness about its healthy uses, possibly about self-managing digital access and study how it can affect stress and, in result, quality of work (increased human error, lack of attention, etc.). AI systems can help understand what is happening and provide the guidance to help steer into the right direction, for increased economic and personal wellbeing.

## 2.2 One Step Further with AI at the Workplace

Current international research and innovation actions (e. g. Horizon 2020 projects SmartWork, WorkingAGE, SustAGE or CO-ADAPT) are focusing on the development and validation of technology-driven interventions, for supporting older employees in their everyday tasks, monitor their health status and enable them to remain longer in the active workforce, while allowing companies to efficiently exploit their long experience and collected knowledge. The SmartWork project (www.smartworkproject.eu/) builds a worker-centric AI system for work ability sustainability of older office workers, by integrating unobtrusive sensing and modeling of the worker state with a suite of novel services for context and worker-aware adaptive work support [10]. On a parallel line, the SustAGE project, develops a multi-modal person-centred IoT platform,

Carina Dantas, Karolina Mackiewicz, Valentina Tageo, Giulio Jacucci, Diana Guardado, Sofia Ortet, Iraklis Varlamis, Michail Maniadakis, Eva de Lera, João Quintas, Otilia Kocsis & Charalampos Vassiliou

which integrates with the daily activities of ageing employees, both at work and outside. The system timely provides its users with personalised recommendations that jointly increase safety, well-being, and productivity of harbor and factory workers. Moreover, in the case of a Covid19-incident for one of the workers, an analysis of all workers positions is carried out in order to identify colleagues who have been close enough to the incident and, therefore, enter a higher control priority. AI can also facilitate the access to knowledge in terms of information and of colleagues, as it is the case in the CO-ADAPT project (https://coadapt-project.eu/) that, besides a personal health assistant, has developed AI (EntityBot) capable to learn in office and knowledge work what tasks the workers in engaged, along with a model of what information entities such as people, documents, applications are connected to it, being then capable to recommend contextually relevant actionable entities that improve the efficiency and quality of work.

Digital tools and services deployed in these projects can pave the way and act as indicative paradigms to similar commercial initiatives on the management of older, as well as younger employees, working on a remote basis. These are AI based tools that bring to the discussion supplementary added-value, but also additional and diverse challenges.

There are enormous benefits of applying AI-based solutions to monitor workers' health and prevent accidents or, currently, COVID-19 infections, and those benefits are reported with enormous potential. According to the recent Deloitte and MedTech Europe report [11], implementing AI in European healthcare systems could save up 380,000 to 403,000 lives annually or €170.9 to 212.4 billion per year.

Moreover, as the paper on "The role of AI technologies in working through COVID-19 and its aftermath" [12] published by Horizon 2020 project SmartWork presented, the AI solutions developed for different scenarios might be particularly useful in the era of pandemic and in the longer future. This non-scientific paper on the role and contribution of the digital solutions and systems to the COVID19 implications in the work environments, gathered all projects funded under the same call to join efforts, reflect and share about the COVID-19 implications to the work environments, now that teleworking turned into a main instrument and necessity for the whole society; understand how the digital solutions and systems could be developed, adapted, optimized or applied to better respond to the pandemic context challenges.

From the different contributions, some similarities can be highlighted: the desire to leverage the existing knowledge and rapidly respond to the challenges of this new (even if hopefully temporary) era; the understanding of the challenges ahead; and the will to overcome them collectively. With the help of technology, employees who have successfully passed through COVID-19 or another high-risk virus, can communicate their experiences in chat-rooms or other social media to alert their colleagues about the possible dangers and make them pay more attention to the prevention procedures.

However, there are also enormous risks of misuse (if not even abuse). Those vary from privacy concerns, gender, disabilities or other discrimination/racial prejudices, or the basic bias based on the poor quality of the data, data collection with inadequate tools, or even the imbalance in power between employer and employees. They should and can be limited in order to fully allow stakeholders to benefit from the opportunities AI solutions may bring. In order to gain the best insights, realistic and massive data collection practices are necessary.

### 2.3 Paving The Way Further

The report on occupational safety and health for EU-OSHA [13] reminds us about all kinds of risks, and that people would prefer AI in the workplace as an on-demand helper rather than as a manager, co-worker or proactive assistant. But, if applied properly, workers believe that AI could improve safety, help reduce mistakes and limit routine work.

The "if applied properly" mentioned earlier is the key differentiator. With the technologies based on personal data, and many of them being so sensitive, such as health related data of various kinds, there is the need to ensure that the solutions are safe, secure, follow the legal standards

Carina Dantas, Karolina Mackiewicz, Valentina Tageo, Giulio Jacucci, Diana Guardado, Sofia Ortet, Iraklis Varlamis, Michail Maniadakis, Eva de Lera, João Quintas, Otilia Kocsis & Charalampos Vassiliou

when it comes to privacy, and at the same time put the individual's control over their data, and their wellbeing in the centre. Clear and transparent communication is necessary to ensure fairness and ethical practices.

Health and productivity go hand in hand and AI can bring benefits for health, such as promoting behaviour change, as well as safety, as long as there is a clear division on the information that is made available to workers and employers.

It is a huge opportunity for AI to understand who is impacted and how, do it at a very early stage and develop all the needed preventive and mitigative actions that may solve some of the challenges, namely sensitive uses, clear accountability, risks to health and safety, potential denial of services.

Four main areas of challenges are thus clustered and highlighted:
- Legal and regulatory - customising solutions at the country level, as well as in different domains, such as data quality, issues of ownership, privacy, ethics and overall data governance.;
- Technical – this includes how data fragmentation can be overcome, storage, access, use, and how to progress on interoperability, not only data quality;
- Social - including workforce accepting and trusting the apps, improving the working environment;
- Education – on the topic of AI for governments, employers and employees.

The societal issues are probably the most challenging hurdles for a wider use of AI in many disciplines, since, as it has been broadly accepted by experts [14]), it is not about the technology but rather about how it is used and governed. Consequently, solving all societal issues is recognised as the crucial issue related to the implementation of AI for health and wellbeing in the workplace.

Complexity is one of the variables that implies societal challenges, mainly connected to the fact that people do not understand what happens to their data nor the benefits – due to the lack of clear and transparent information. Trust in the digital tools and the use of data is the main challenge to developers and there is the need for broader user validation in AI to increase trust.

One additional angle that requires discussion is to understand how big is the risk that AI solutions can endanger jobs if they massively analyse productivity? And how can misuse by employers or other authorities be prevented? However, there are technical ways in which good governance can call for anonymization of data in a way that it can still contribute to providing insights that help improve existing processes, without pointing at specific people. Hence the insistence in proper data privacy policies and governance in place. Employees need to understand the company's practices in their decisions.

Departing from the book "Architects of Intelligence: The truth about AI from the people building it" [15] it is also worth to point out that it is the first time in history that more jobs are erased by a new technology than the ones created, which means that this a real discussion to hold. Automation may provoke less jobs, a shift in job types and creations, a new universe of jobs has grown as a result, while others are rapidly fading, and measures to prevent or address this in the societal area are required. Is a minimum global income to protect people that may lose jobs one solution to discuss? Will it be accepted that the State will provide money for those who do nothing? How is the digital transition managed?

On the crossover between the human and technological factors, an important question for apps and systems to work correctly is eliciting the correct rules and defining the right values. AI technologies don't all work the same, or face the same challenges. For example, machine learning algorithms deal with a huge amount of data, which can be anonymized, but still rely on data quality, whereas chatbots can focus on individual and personal data and encode rules and knowledge in order to interpret human-provided input. Consequently, the performance of AI

Carina Dantas, Karolina Mackiewicz, Valentina Tageo, Giulio Jacucci, Diana Guardado, Sofia Ortet, Iraklis Varlamis, Michail Maniadakis, Eva de Lera, João Quintas, Otilia Kocsis & Charalampos Vassiliou

solutions strongly depends on the data provided and governance, as well as on the models developed, but in any case, there is still a long way to go. Privacy is, of course, a relevant issue, but to enable trust on AI, the collection instruments and the transmission channels are really important. Recent developments like blockchain technology may provide solutions with respect to data privacy and protection in order to improve citizens' trust on AI.

This also highlights the need for more clarity in what concerns the EU Health Data Space [16] and the different uses of data, anonymised or personal. Fair data and standards bring also very important initiatives in Europe to enable progress in this area. Data governance models are necessary to help guide the transition into AI-friendly ecosystems and organizations.

## 2.4    Discussion
When applied to the work environment, AI poses several possibilities related to its ability, not only to learn and predict human behaviour, but also to develop its own value and 'morals'. Under the actual pandemic context, which seems to be here to stay for a while, the AI potential allows us to daydream on how, at the workplace, it could be used to promote and ensure the compliance of contingency measures and safety procedures, by workers, employers and even clients within the organizations, as well as to maintain or increase the workability sustainability of the oldest ones in the active.

On the other hand, at home or within other personal contexts (e.g. outside, while shopping, or in the gym), AI solutions could be conceived and personalized to ensure physical distancing, to detect risk behaviours, or to avoid unnecessary travelling or face to face interaction, thus avoiding the spread of Coronavirus. In addition, AI could also contribute to prevent older workers (at greater risk towards this threat) from becoming dangerously isolated, lacking social interaction or family support, by enhancing their linking to the people and services in the outside world.

However, apart from the good performance of an AI solution, humans need to 'see inside' the black box and understand how and why. This raises the need for explainable AI. People need to be at the center of the systems to be able to progress, which also leads to accountability – what are the systems doing? Who owns this information? This needs to be fully understood to ensure that systems and apps are used and uptake. Nevertheless, misuse can only be established through a regulatory framework at national level through the definition of an ethical framework.

In practice, this means:
  •   Developing ethical standards and policy frameworks to build trust and foster the adoption of AI in healthcare in conjunction with the working environment,
  •   Securing access to the high-quality data by building data policies and infrastructure to foster access, and interoperability of the harmonised data;
  •   Respecting the employees' rights to privacy and confidentiality by making sure the data is collected and managed properly, and used meaningfully and in compliance with their fundamental human rights and informed consent;
  •   Improving explainability and accountability, as well as digital literacy among all stakeholders involved, from the decision makers to the employers to the employees.

Although there is already extensive publication in this field, most of it is either devoted to the technical aspects or only to the social or ethical ones. An integrated approach to AI is still needed and this collaborative and multidisciplinary understanding of the challenges and benefits that is at the heart of this article, intends to pave the way beyond a siloed perspective to call for a concerted policy agenda and implementation strategy of AI in the workplace.

## 2.5    Conclusions
The pandemic crisis extensively affected the way companies operate, transforming them from premise-centered to remote (work from home)-centered business. The use of technology and available digital tools enabled this transition, without the need, surprisingly, to seriously disrupt the ongoing operations.

Carina Dantas, Karolina Mackiewicz, Valentina Tageo, Giulio Jacucci, Diana Guardado, Sofia Ortet, Iraklis Varlamis, Michail Maniadakis, Eva de Lera, João Quintas, Otilia Kocsis & Charalampos Vassiliou

However, the success of this digital transformation significantly depends on how effectively challenges related with the established procedures, the organisation's nature and the safety of the employees are tackled and further elaboration will be required to adjust to these unprecedented working conditions.

Even though COVID served as an accelerator for the way we worked, moving from the traditional office space to a hybrid or online model, the increased adoption of the internet and communication platforms had already introduced and pushed for a tendency to working remotely. As the work environment transfers from a one-context to a multi-context (including the street), different needs arise to help balance work and personal life, without reducing quality, but the opposite, aiming at increased work efficiency and satisfaction, and greater quality of life.

Moreover, the existing digital gap, a huge number of people with problems accessing and using technologies (15% of the world's population according to the WHO [17]), suffered tremendously with this dramatic shift, because homes and other contexts where not prepared for working remotely. For example, people with disabilities had special assistive technologies in the workplace but not in their home environment. Designing inclusive solutions, including AI, in all fields of life and contexts, is now more critical than ever.  The collection of data to help generate insights, in the form of information, can greatly help save time, increase efficacy and efficiency, and result in a smoother transition to this new lifestyle.

This is a call to action for a holistic and fully inclusive approach to AI, that goes beyond technology and includes ethical implementation, user-centricity, cross-sectoral policies and limits the risks to capitalise on the benefits of the new technologies.

However, the post-pandemic economic crisis and the social changes that are emerging from this period will easily create the opportunity to remove these priorities from the political agenda and funding plans and that, unless action is taken, can lead to a setback of more than five years in implementing innovation and quality of life for citizens.

A huge ethical challenge to be faced will be to redefine the balance between digital tools and human presence. If this was already somewhat stable in public opinion, the emergency period polarized opinions once again and this may actually be a threat to the broader adoption of A.I. tools for an increased workability.

It will take an enormous sensitivity and a great social conscience to evolve in the right direction and not lose focus - all political, economic, and social measures must have the ultimate goal of people's wellbeing and the promotion of common good.

## 3.  ACKNOWLEDGEMENTS

## 4.  REFERENCES
[1] McKinsey & Company, Unlocking success in digital transformations https://www.mckinsey.com/business-functions/organization/our-insights/unlocking-success-in-digital-transformations#, Oct. 29, 2019 [May. 05, 2021].

[2] R. Vadivel & Dr K. Baskaran. "Three Dimensional Database: Artificial Intelligence to eCommerce Web service Agents", in International Journal of Computer Science and

Carina Dantas, Karolina Mackiewicz, Valentina Tageo, Giulio Jacucci, Diana Guardado, Sofia Ortet, Iraklis Varlamis, Michail Maniadakis, Eva de Lera, João Quintas, Otilia Kocsis & Charalampos Vassiliou

Security (IJCSS), Volume (5) : Issue (2) : 2011 181 http://www.cscjournals.org/manuscript/Journals/IJCSS/Volume5/Issue2/IJCSS-422.pdf.

[3]     A. Silva-C et al. "The attitude of managers toward telework, why is it so difficult to adopt it in organizations?", in Technology in Society, Volume 59, 2019.

[4]     A. Belzunegui-Eraso, A. Erro-Garcés. "Teleworking in the Context of the Covid-19 Crisis. Sustainability". 2020; 12(9):3662.

[5]     M. Fana, et al. "Working Paper Telework, work organisation and job quality during the COVID-19 crisis: a qualitative study" JRC Working Papers Series on Labour, Education and Technology, No. 2020/11, European Commission, Joint Research Centre (JRC), Seville

[6]     M. de Haas, R. Faber, M. Hamersma, "How COVID-19 and the Dutch 'intelligent lockdown' change activities, work and travel behaviour: Evidence from longitudinal data in the Netherlands", Transportation Research Interdisciplinary Perspectives, Volume 6, 2020, 100150, ISSN 2590-1982.

[7]     S. Shah, D. Nogueras, H. van Woerden, V. Kiparoglou. "The COVID-19 Pandemic: A Pandemic of Lockdown Loneliness and the Role of Digital Technology", J Med Internet Res 2020;22(11):e22287.

[8]     A. Ancillo, M. Núñez & S. Gavrila. "Workplace change within the COVID-19 context: a grounded theory approach", Economic Research-Ekonomska Istraživanja, 2020.

[9]     D. Reuschke, A. Felstead. "Changing workplace geographies in the COVID-19 crisis", Dialogues in Human Geography, 2020, Vol. 10(2) 208–212.

[10]   O. Kocksis, et al. "SmartWork: Designing a Smart Age-Friendly Living and Working Environment for Office Workers" in PETRA '19: Proceedings of the 12th ACM International Conference on PErvasive Technologies Related to Assistive Environments, June 2019, pp 435–441.

[11]   Deloitte. "The socio-economic impact of AI in healthcare", https://www.medtecheurope.org/wp-content/uploads/2020/10/mte-ai_impact-in-healthcare_oct2020_report.pdf, 2020 [Mar. 29, 2021].

[12]   SmartWork project. "Concerted Paper Between Partners and Projects, The role of AI technologies in working through COVID-19 and its aftermath", https://ec.europa.eu/info/funding-tenders/opportunities/portal/screen/opportunities/horizon-results-platform/26019;resultId=26019;campaigns=1, 2020 [Mar. 29, 2021].

[13]   P.V.Moore "OSH and the Future of Work: Benefits and Risks of Artificial Intelligence Tools in Workplaces", https://osha.europa.eu/en/publications/osh-and-future-work-benefits-and-risks-artificial-intelligence-tools-workplaces/view, Jul. 5, 2019 [Mar. 29, 2021].

[14]   K. Mackiewicz and C. Dantas. "Endless Benefits Vs. Huge Risks: AI And Health At Workplace", https://ictandhealth.com/news/endless-benefits-vs-huge-risks-ai-and-health-at-workplace/, Dec. 11, 2020 [Mar. 29, 2021].

[15]   M. Ford. "Architects of Intelligence: The truth about AI from the people building it". Packt Publishing, Nov. 23, 2018.

[16]   European Commission. "European Health Data Space", https://ec.europa.eu/health/ehealth/dataspace_en, 2020, [Mar. 29, 2021].

[17]   World Health Organisation. "Disability and health", https://www.who.int/news-room/fact-sheets/detail/disability-and-health, 2020, [Mar. 29, 2021].

Sundar Krishnan, Ashar Neyaz & Qingzhong Liu

# IoT Network Attack Detection using Supervised Machine Learning

**Sundar Krishnan**                                                      *skrishnan@shsu.edu*
*Department of Computer Science*
*Sam Houston State University*
*Huntsville, TX, USA*

**Ashar Neyaz**                                                      *ashar.neyaz@shsu.edu*
*Department of Computer Science*
*Sam Houston State University*
*Huntsville, TX, USA*

**Qingzhong Liu**                                                      *liu@shsu.edu*
*Department of Computer Science*
*Sam Houston State University*
*Huntsville, TX, USA*

**Abstract**

The use of supervised learning algorithms to detect malicious traffic can be valuable in designing intrusion detection systems and ascertaining security risks. The Internet of things (IoT) refers to the billions of physical, electronic devices around the world that are often connected over the Internet. The growth of IoT systems comes at the risk of network attacks such as denial of service (DoS) and spoofing. In this research, we perform various supervised feature selection methods and employ three classifiers on IoT network data. The classifiers predict with high accuracy if the network traffic against the IoT device was malicious or benign. We compare the feature selection methods to arrive at the best that can be used for network intrusion prediction.

**Keywords:** Supervised Learning, Network Attack Detection, IoT, Network Forensics, Network Security.

## 1. INTRODUCTION
Network traffic has seen unprecedented growth in the last decades. With growing volumes of Internet-connected devices, cheaper cloud storage, growing smartphone technology, decreasing device and network hardware costs, and the advent of 5G technology, it is predicted that by 2023, there will be 3X more networked devices on earth than humans. A Cisco Annual Internet Report Forecasts 5G to support more than 10% of Global Mobile Connections by 2023 [1], [2]. This growth in network traffic and Internet-connected devices has resulted in an increase of malicious attacks over the network that can sometimes be difficult to detect. A network attack is a type of cyber-attack in which the attacker attempts to gain unauthorized access into a computer network or an Internet-connected device for malicious purposes or reconnaissance. Cyber-attacks rank as the fastest growing crime in the U.S., causing catastrophic business disruption. Globally, cybercrime damages are expected to reach US $10.5 trillion annually by 2025. NIST defines a cyber-attack (breach) as, "An attack, via cyberspace, targeting an enterprise's use of cyberspace for the purpose of disrupting, disabling, destroying, or maliciously controlling a computing environment/infrastructure; or destroying the integrity of the data or stealing controlled information" [3]. Over the years, Cybercrime has moved on from targeting and harming people, computers, networks, and smartphones - to cars, power-grids, smart devices, and anything that can connect back to the Internet.

Sundar Krishnan, Ashar Neyaz & Qingzhong Liu

The Internet of Things (IoT) has come a long way since the 80s when early IoT designers (students) at Carnegie Melon University installed micro-switches inside of a Coca-Cola vending machine to remotely check on the temperature and availability of their favorite beverages [4]. IoT devices and technology have gone mainstream these days, with IoT devices remotely controlling our home speakers, smart elevators, cars, household appliances, power plants, security cameras, baby cams, smart buildings, medical devices, freight, etc. These devices connect back to the Internet via traditional copper wires, fiber, and telecom technology for remote control functionality, thereby making them game for malicious actors using the Internet. IoT devices are often shipped to users with minimal logon security, operating system vulnerabilities, and overall poor security design. This can be mostly attributed to keeping costs down, ease of use for the user, and inadequate security foresight by the manufacturers. Consequently, the attack surface of IoT devices has greatly grown, triggering security and privacy concerns. The infamous Mirai botnet [5] self-replicated by seeking out hundreds of thousands of home routers with weak or non-existent passwords. The roll-out of the 5G mobile networks may further embolden IoT cyber attackers due to the advantage of high bandwidth, ultra-low latency, and fundamentally new networking capabilities of 5G technology [6].

IoT tangibly solves many business problems across industries such as healthcare, smart cities, building management, utilities, transportation, and manufacturing. About 30% of devices on enterprise networks today are network-connected IoT devices [7], making them potential targets over a network. Unlike traditional IT assets like servers and endpoints, these network-connected devices may not be well maintained and documented by IT teams. Such assets can easily be missed from an organization's proactive security monitoring apparatus. IoT devices are also found in home networks and may not have adequate security controls or infrastructure to protect them. With more and more diverse types of IoT devices continuing to connect to the network, there can be a dramatic broadening of the attack surface. All it takes for a successful intrusion is the diminished integrity of a weak asset on the network.

Predictive capabilities are incredibly beneficial in any industrial setting, especially in thwarting cyber-attacks. Machine learning helps solve tasks (such as regression, clustering, classification, dimensionality reduction, etc.) using an approach/method based on available data. A popular area of machine learning application in cybersecurity is helping businesses detect malicious activity faster and stop attacks before they get started. Cybersecurity should be implemented in layers against any asset. It must be noted that machine learning alone will never be a silver bullet for cybersecurity, but when coupled with other controls, it can improve intrusion detection. While extensive research has been undertaken to predict/detect network attacks on common Information Technology assets, little research has been conducted towards IoT device attacks. In this research, we apply machine learning approaches towards IoT attack detection using the IoTID20 dataset [8] that was built on the network traffic of botnet attacks [9] against IoT devices. Three feature selection models are chosen, and the prediction of an attack based on supervised learning is presented by applying three classifiers against each feature selection model.

## 2. BACKGROUND

An intrusion detection system (IDS) is a hardware device or software application that monitors a network or host for malicious activity or policy violations [10]. While IDS alerts on intrusions, Intrusion Prevention Systems (IPS) can respond to detected intrusion upon discovery. Intrusion detection using both supervised learning and unsupervised learning has been primarily researched. Using unsupervised machine learning to understand better network attacks has been widely researched. Kumar et al. [11] evaluated MeanShift algorithm to detect network incursion against the KDD99 network traffic dataset. The authors concluded that the MeanShift could detect an attack in the network dataset. However, the algorithm could not detect Remote to Local (R2L) and User to Root (U2R) attacks. Serra et al. Mukherjee et al. [12] proposed ClusterGAN as a new medium for adaptive clustering using Generative Adversarial Networks (GANs). Choi et al. developed a network intrusion detection system (NIDS) using an unsupervised learning algorithm against unlabeled data. The high accuracy of the experiment results provided a recommendation

for developing network intrusion detection systems. False attack detection can be challenging to detect. Sakhini et al. [14] evaluated SVM (Support Vector Machine), KNN (K-Nearest-Neighbors), and ANN (Artificial Neural Network) to detect FDI (False Data Injection) attacks. Their experiment results showed that KNN and SVM were more accurate than ANN. Supervised learning is the machine learning task of learning a function that maps an input to an output based on examples (labeled data) of such input-output pairs. Balkanli et al. [17] detected network intrusion with 99% accuracy against 20% of backscatter darknet traffic by employing two opensource network intrusion detection systems (NIDS) and two supervised machine learning techniques on backscatter darknet traffic. Morfino et al. [18] evaluated the performance of various supervised machine learning algorithms in recognizing cyberattacks, specifically, the SYN-DOS attacks on IoT systems by differentiating them in terms of application performances and also in training/application times. Their Apache Spark algorithm yielded an accuracy of greater than 99%, whereas Random Forest achieved an accuracy of 1%. A simple type of attack against IoT devices is Denial-of-Service (DoS). The IoT device receives bursts of surplus network traffic rendering it unusable or overtaxing IoT hardware and underlying infrastructure. Hodo et al. [19] used Artificial Neural Network (ANN) to detect Denial-of-Service (DoS) of Distributed Denial-of-Service (DDoS) attacks with a 99.4% accuracy in attack detection. Loannou et al. [20] put forward the use of Support Vector Machine (SVM) learning model for detecting deviation within the Internet of Things. The proposed SVM model achieved up to 100% accuracy when evaluated against the unknown data taken from the corresponding network topology with proper training. The model also achieved an 81% accuracy when used under an unknown topology. Often IoT devices are wireless and configured to routers with poor security settings [21]. Grimaldi et al. [22] leveraged supervised machine learning techniques in real-time to identify and detect wireless traffic interference, thereby allowing for isolation and extraction of standard-specific traffic. Anthi et al. [23] presented a three-layer intrusion detection system (IDS) that used a supervised machine learning approach to detect a variety of popular network-based attacks on IoT networks. The proposed system's three core functions' performance resulted in an F-measure of 96.2%, 90.0%, and 98.0%, respectively. This demonstrated that the proposed system could automatically distinguish IoT devices on the network and detect attack types against devices on the network. Artificial neural networks and deep learning approaches can also be used to detect network intrusions. Caron et al. [16] proposed a scalable clustering method called DeepCluster for unsupervised learning of convolutional neural networks or convnets against the ImageNet and YFCC100M datasets. Their results obtained were better than other state-of-the-art approaches by a significant margin. There are limitations to using machine learning to identify network attacks. Xiao et al. [15] examined attack models and IoT security solutions based on machine learning techniques. They concluded that supervised and unsupervised learning sometimes fails to detect the attacks due to oversampling, insufficient training data, and bad feature extraction. In this research, we leverage supervised learning to predict normal and malicious/abnormal network traffic using the IoTID20 dataset [8]. Ullah et al. [8] proposed this dataset, namely IoTID20 that was generated from Botnet traffic against IoT devices [9]. Ullah et al. [8] also utilized this dataset to propose a detection classification methodology. In this article, we choose a different approach compared to Ullah et al. [8] when selecting features and classifiers. We then evaluate these various feature selection approaches against classifier accuracy.

## 3. FEATURE SELECTION IN MACHINE LEARNING

Machine learning is a branch of computational algorithms designed to emulate human intelligence by learning from the surrounding environment [24]. Machine learning (ML) and Artificial Intelligence (AI) have become dominant problem-solving techniques in many areas of research and industry in the last decade. ML and AI are not the same. While Artificial intelligence is about problem-solving, reasoning, and learning in general; Machine learning is specifically about learning—learning from examples, from definitions, from being told, and from behavior [25]. While working with ML, we typically use datasets (like a database table or an Excel spreadsheet) that contain data for the experiment arranged in columns (features). Each feature, or column, represents a measurable piece of data that can be used for analysis. The below discussion is about a few feature engineering (selection) techniques and supervised learning algorithms

employed in our research experiments. Often in a dataset, the given set of features in their raw form does not provide enough, or the most optimal, information to train a Machine Learning model. It may be beneficial to remove unnecessary or conflicting features in some instances, which is known as feature selection or feature engineering. Feature selection is a critical and effective approach to ignoring or retaining certain features on a dataset that do not contribute statistically significantly towards the predicted outcome. Thus, only the most significant subset of features are retained in a model while removing these irrelevant, redundant, and noisy features.

### 3.1 Filter Methods
Filter methods select features from a dataset independently by relying on features' characteristics, which is often the first step before applying machine learning algorithms. Basic and intuitive filter methods help remove Constant features, Quasi- Constant features, and Duplicated features. A dataset can also include correlated features wherein highly correlated features provide redundant information regarding the target. In such cases, removing one of the two highly correlated features can reduce the dimensionality and noise.

### 3.2 Sequential Forward Processing
Sequential Forward Processing (or forward feature selection) is a wrapper method that iterates through the set of features while evaluating them using a machine learning algorithm. A preset criterion (k features) is selected, which is the maximum number of features to be reached when starting from zero. The initial starting step is to evaluate all features individually and then select the one that results in the best performance [26]. In the second iteration, we test all possible combinations of the selected feature with the remaining features and retain the pair that produces the best algorithmic performance. Subsequent iterations continue by adding one feature at a time in each iteration until the preset criteria is reached.

### 3.3 Sequential Backward Processing
Sequential Backward Processing (or backward feature selection) is a wrapper method that iterates through the set of features while evaluating them using a machine learning algorithm. A preset criterion (k features) is selected, which is the maximum number of features to be reached when starting from zero. The initial starting step is to consider all the features in the dataset, followed by a performance evaluation of the algorithm [26]. Similar iterations follow by removing one feature (least significant) at a time producing the best performing algorithm using an evaluation metric. Iterations continue removing feature after feature until the preset criteria is reached.

### 3.4 Recursive Feature Elimination
Recursive Feature Elimination (RFE) is a feature selection method that fits a model (e.g., linear regression or SVM) and removes the weakest feature (or features) until the specified number of features are reached [27]. RFE requires a specified number of features to keep while eliminating dependencies and collinearity that may exist in the model.

## 4.  SUPERVISED LEARNING
Supervised learning in machine learning and artificial intelligence refers to systems and algorithms that determine a predictive model using labeled data points with known outcomes. The model is learned by training through learning algorithms such as linear regression, random forests, or neural networks. As input data is fed into the model, it adjusts its weights through a reinforcement learning process, ensuring that the model has been fitted appropriately [28]. Supervised learning is often used to create machine learning models for Regression and Classification types of problems. A statistical approach known as regression analysis can be implemented to establish a possible relationship between different variables. Regression analysis consists of a set of machine learning methods that allow predicting a continuous outcome variable (y) based on the value of one or multiple predictor variables (x).

### 4.1 Random Forest

Random Forest (RF) is based on decision trees and is one of the many machine learning algorithms used for supervised learning. There are two main ways for combining the outputs of multiple decision trees into a random forest, 1. Bagging (Bootstrap aggregation) used in Random Forests and 2. Boosting (used in Gradient Boosting Machines). Figure: 1 depicts a random forest. Random Forest implementations are available in many machine learning libraries for R and Python, like Caret (R) [30], Scikit-learn (Python sklearn.ensemble.RandomForestRegressor) [31], and H2O (R and Python) [32].
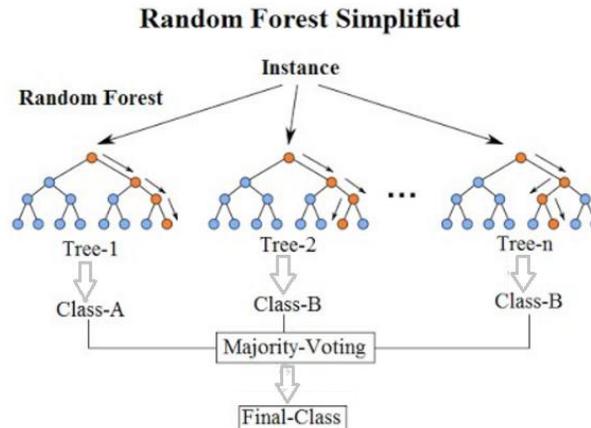


**FIGURE 1:** A diagram of a random decision forest [29].

### 4.2 Support Vector Classifier (SVC)

The main task of the algorithm is to find the most correct line, or hyperplane, which divides data into two classes. An SVC is an algorithm that receives input data and returns such a dividing line. In python sklearn library [31], the implementation of SVC is based on libsvm. The objective of a Linear SVC (Support Vector Classifier) is to fit the data and return a "best fit" hyperplane that divides or categorizes the data.

### 4.3 eXtreme Gradient Boosting (XGBoost)

XGBoost implements machine learning algorithms under the gradient boosting framework and is an optimized end-to-end tree boosting library designed to be highly efficient, flexible, and portable [33]. The XGBoost library implements the gradient boosting decision tree algorithm. Figure: 2 depicts the evolution of XGBoost. Generally, XGBoost is fast when compared to other implementations of gradient boosting [35].
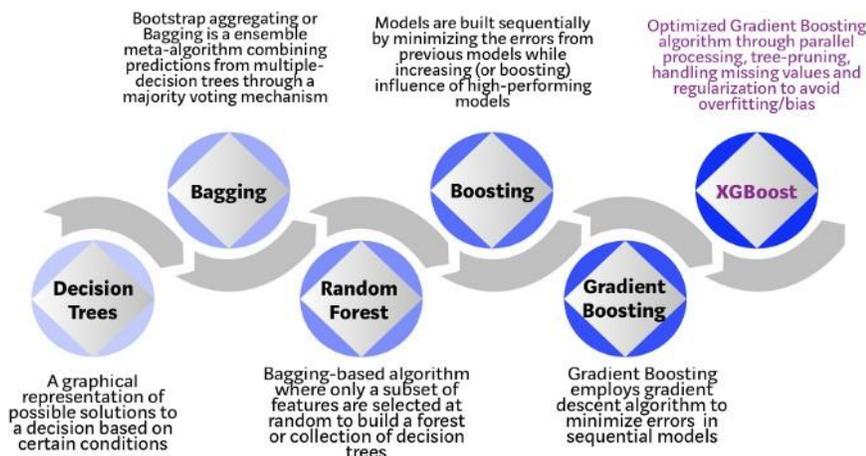


**FIGURE 2:** Evolution of XGBoost Algorithm from Decision Trees [34].

## 5. EXPERIMENTS

The IoT dataset in .csv format obtained by Ullah et al. [8] was used for this experiment. The dataset contains network traffic processed from packet captures [9] on two smart home devices wherein attacks on these IoT devices were captured over the wireless network. We decided to focus on the month of May for its ease of use and as it contained all the necessary network attack categories. Python scripts were used for parsing, data preparation, and logistic regression. Results were then documented for analysis. Figure: 3 outlines the workflow of our research.

### 5.1 Dataset Preparation

Before logistic regression analysis could be performed, few data preparation steps were undertaken below to pre-process and transform the raw data into the necessary data structure to carry out the analysis. The timestamp feature was first formatted for a timestamp format. The dataset was filtered for May/2019 network traffic data using the Timestamp feature. We decided to ignore the features FlowID, Category, and Sub-Category. The feature Label was encoded for Normal =1 and Anomaly=2. The Src IP and Dest IP features were each encoded as 1, 2, 3, and 4 depending on the network class of the IP address values (class A=1, B=2, C=3 and D=4). The Timestamp feature was transformed into Date and Time features (24-hr format). Data rows with invalid Dst IP (0.0.0.x) were ignored. Table 1 shows the features at the end of this step. We used pre-processing techniques such as dropping features that are Constants, Quasi-Constants, and Duplicates. A Pearson's correlation of 0.8 was then applied to further select features. Correlated features degrade the detection capability of a machine learning algorithm, and thus, highly correlated features were ignored from the IoTID20 dataset. Table 2 shows a list of features that were dropped from the IoTID20 dataset at each pre-processing stage and the final set of features to retain at the end of pre-processing.
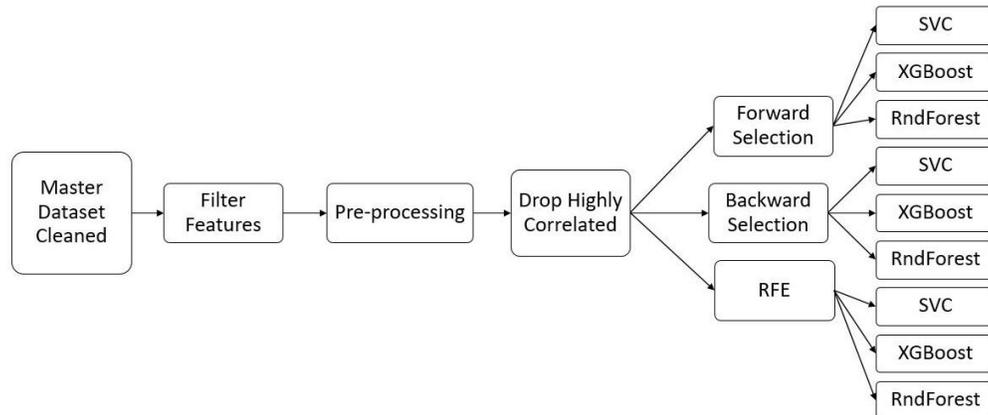


**FIGURE 3:** Experiment workflow.

### 5.2 Feature Selection and Logistic Regression

After pre-processing, a separate dataset with 33 features as in Table 3 was created for the experiment. This dataset was then used for the below experiments.

a) Applied Sequential Backward Processing for feature selection and obtained eight features for logistic regression. Table 2 shows the features obtained after applying Sequential Backward Processing. The dataset was split into train/test (80/20) and perform logistic regression using Random Forest Classifier, SVC, and XGBoost. Results were documented for analysis.

b) Applied Sequential Forward Processing for feature selection and obtained eight features for logistic regression. Table 2 shows the features obtained after applying Sequential Forward Processing. The dataset was split into train/test (80/20) and perform logistic regression using Random Forest Classifier, SVC, and XGBoost. Results were documented for analysis.

c) Applied RFE processing for feature selection and obtained eight logistic regression features. Table 2 shows the features obtained after applying RFE. The dataset was split into train/test (80/20) and perform logistic regression using Random Forest Classifier, SVC, and XGBoost. Results were documented for analysis.

| Dataset | Features |
|---------|----------|
| Original dataset | Src IP,Src IP Cl (network class of Src IP), Src Port, Dst IP, Dst IP CL (network class of Dst IP), Dst Port, Protocol, Timestamp DT (split date value of Timestamp), Timestamp 24HR TIME (split time value of Timestamp), Flow Duration, Tot Fwd Pkts, Tot Bwd Pkts, TotLen Fwd Pkts, TotLen Bwd Pkts, Fwd Pkt Len Max, Fwd Pkt Len Min, Fwd Pkt Len Mean, Fwd Pkt Len Std, Bwd Pkt Len Max, Bwd Pkt Len Min, Bwd Pkt Len Mean, Bwd Pkt Len Std, Flow Byts s, Flow Pks s, Flow IAT Mean, Flow IAT Std, Flow IAT Max, Flow IAT Min, Fwd IAT Tot, Fwd IAT Mean, Bwd IAT Mean, Fwd IAT Max, Fwd IAT Min, Bwd IAT Tot, Bwd IAT Mean.1, Bwd IAT Std, Bwd IAT Max, Bwd IAT Min, Fwd PSH Flags, Bwd PSH Flags, Fwd URG Flags, Bwd URG Flags, Fwd Header Len, Bwd Header Len, Fwd Pkts s, Bwd Pks s, Pkt Len Min, Pkt Len Max, Pkt Len Mean, Pkt Len Std, Pkt Len Var, FIN Flag Cnt, SYN Flag Cnt, RST Flag Cnt, PSH Flag Cnt, ACK Flag Cnt, URG Flag Cnt, CWE Flag Count, ECE Flag Cnt, Down Up Ratio, Pkt Size Avg, Fwd Seg Size Avg, Bwd Seg Size Avg, Fwd Byts/b Avg, Fwd Pkts b Avg, Fwd Blk Rate Avg, Bwd Byts b Avg, Bwd Pkts b Avg, Bwd Blk Rate Avg, Subflow Fwd Pkts, Subflow Fwd Byts, Subflow Bwd Pkts, Subflow Bwd Byts, Init Fwd Win Byts, Init Bwd Win Byts, Fwd Act Data Pkts, Fwd Seg Size Min, Active Mean, Active Std, Active Max, Active Min, Idle Mean, Idle Std, Idle Max, Idle Min, Label (converted to binary), Cat, Sub Cat |

**TABLE 1:** Original Dataset after feature cleansing.

| Pre-processing / Filter techniques applied | Features dropped as a result |
|---------------------------------------------|------------------------------|
| Constant features | Fwd PSH Flags, Fwd URG Flags, Fwd Byts/b Avg, Fwd Pkts b Avg, Fwd Blk Rate Avg, Bwd Byts b Avg, Bwd Pkts b Avg, Bwd Blk Rate Avg, Init Fwd Win Byts, Fwd Seg Size Min, Timestamp month |
| Quasi-Constant features | Bwd URG Flags, FIN Flag Cnt, RST Flag Cnt, URG Flag Cnt, CWE Flag Count, ECE Flag Cnt |
| Duplicate Features | PSH Flag Cnt, Fwd Seg Size Avg, Bwd Seg Size Avg, Subflow Fwd Pkts, Subflow Fwd Byts, Subflow Bwd Pkts, Subflow Bwd Byts |
| Correlated features | Bwd Pkt Len Mean, Idle Mean, Idle Std, Bwd IAT Tot, Fwd Act Data Pkts, Active Min, Fwd Header Len, ACK Flag Cnt, Bwd Pks s, Pkt Len Var, Active Max, Flow IAT Min, Timestamp day, Idle Max, Label, Idle Min, Fwd Pkt Len Mean, TotLen Fwd Pkts, Pkt Len Mean, Flow IAT Max, Bwd Pkt Len Min, Pkt Len Max, Pkt Size Avg, Fwd IAT Min, Bwd IAT Max, Fwd IAT Mean, Timestamp hour, Bwd IAT Min, Fwd IAT Max, Flow Byts s, Bwd IAT Mean.1 |
| **Final set of features for experiment** | |
| Src IP Cl, Src Port, Dst IP CL, Dst Port, Protocol, Flow Duration, Tot Fwd Pkts, Tot Bwd Pkts, TotLen Bwd Pkts, Fwd Pkt Len Max, Fwd Pkt Len Min, Fwd Pkt Len Std, Bwd Pkt Len Max, Bwd Pkt Len Std, Flow Pks s, Flow IAT Mean, Flow IAT Std, Fwd IAT Tot, Bwd IAT Mean, Bwd IAT Std, Bwd PSH Flags, Bwd Header Len, Fwd Pkts s, Pkt Len Min, Pkt Len Std, SYN Flag Cnt, Down Up Ratio, Init Bwd Win Byts, Active Mean, Active Std, Label, Timestamp minute, Timestamp second | |

**TABLE 2:** Pre-processing of the IoT20 Dataset.

| Feature Selection Applied | Features Obtained |
|---|---|
| Sequential Backward Processing | 'Pkt Len Var', 'Fwd Header Len', 'TotLen Fwd Pkts', 'Timestamp hour', 'Pkt Len Mean', 'Fwd Pkt Len Mean', 'Bwd Pkt Len Min', 'Timestamp day', 'Label' |
| Sequential Forward Processing | 'Pkt Len Var', 'Timestamp hour', 'Pkt Len Mean', 'Fwd Pkt Len Mean', 'Bwd Pkt Len Min', 'Pkt Size Avg', 'Timestamp day', 'Bwd Pkt Len Mean', 'Label' |
| RFE | 'Flow IAT Max', 'ACK Flag Cnt', 'Idle Max', 'Timestamp hour', 'Fwd Act Data Pkts', 'Flow IAT Min', 'Idle Min', 'Idle Std', 'Label' |

**TABLE 3:** Feature selection techniques applied for Logistic Regression.

## 6. ANALYSIS

After pre-processing, the dataset was put to three feature selection processes to arrive at eight highly ranked features. This was followed by three logistic regression algorithms. The first feature selection method was Sequential Backward Processing. The results of SVC, XGBoost, and Random Forest classification against the eight features from Sequential Backward Processing are shown in Table 6. The ROC curve for the three classifiers is shown in Figure: 5, and the Reliability Curve is shown in Figure 4. The second feature selection method was Sequential Forward Processing. The results of SVC, XGBoost, and Random Forest classification against the eight features from Sequential Forward Processing are shown in Table: IV. The ROC curve for the three classifiers is shown in Figure 7, and the Reliability Curve is shown in Figure 6. The third feature selection method was Recursive Feature Elimination (RFE). The results of SVC, XGBoost, and Random Forest classification against the eight features from RFE are shown in Table: IV. The ROC curve for the three classifiers is shown in Figure 9, and the Reliability Curve is shown in Figure 8. From Table 4, we can conclude that all the three supervised feature selection methods could predict with high accuracy malicious traffic vs. benign traffic. The number of features used (eight) was random, but changes to the number used can impact accuracy scores. The Root Mean Square Error (RMSE) is a valuable metric that tells us how far apart our predicted values are from our observed values in a model. The SVC classifier has larger RMSE values in the three feature selection methods, implying a worse model fits the data. Overall, the use of RFE yielded the best accuracy for the three classifiers.

| Feature Selection Method | Classifier | Accuracy | F1 Score | Recall | RMSE |
|---|---|---|---|---|---|
| Sequential Backward Processing | SVC | 98.20% | 0.98 | 0.97 | 0.134096 |
| | XGBoost | 99.31% | 0.99 | 0.98 | 0.082918 |
| | Random Forest | 99.23% | 0.99 | 0.98 | 0.087708 |
| Sequential Forward Processing | SVC | 98.48% | 0.98 | 0.97 | 0.123455 |
| | XGBoost | 99.30% | 1.00 | 0.98 | 0.083495 |
| | Random Forest | 99.21% | 0.99 | 0.98 | 0.089068 |
| Recursive Feature Elimination | SVC | 98.76% | 0.98 | 0.98 | 0.111159 |
| | XGBoost | 99.79% | 1.00 | 1.00 | 0.04599 |
| | Random Forest | 99.78% | 1.00 | 1.00 | 0.047028 |

**TABLE 4:** Logistic Regression Results.

Sundar Krishnan, Ashar Neyaz & Qingzhong Liu



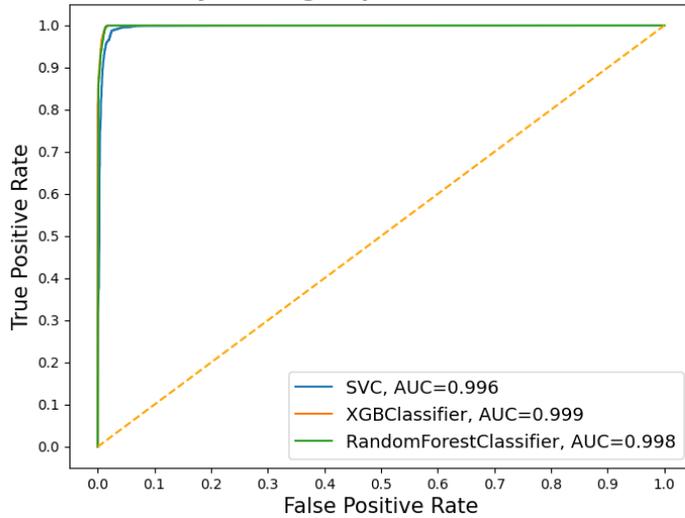**FIGURE 4:** Reliability Curve after applying sequential backward feature selection.



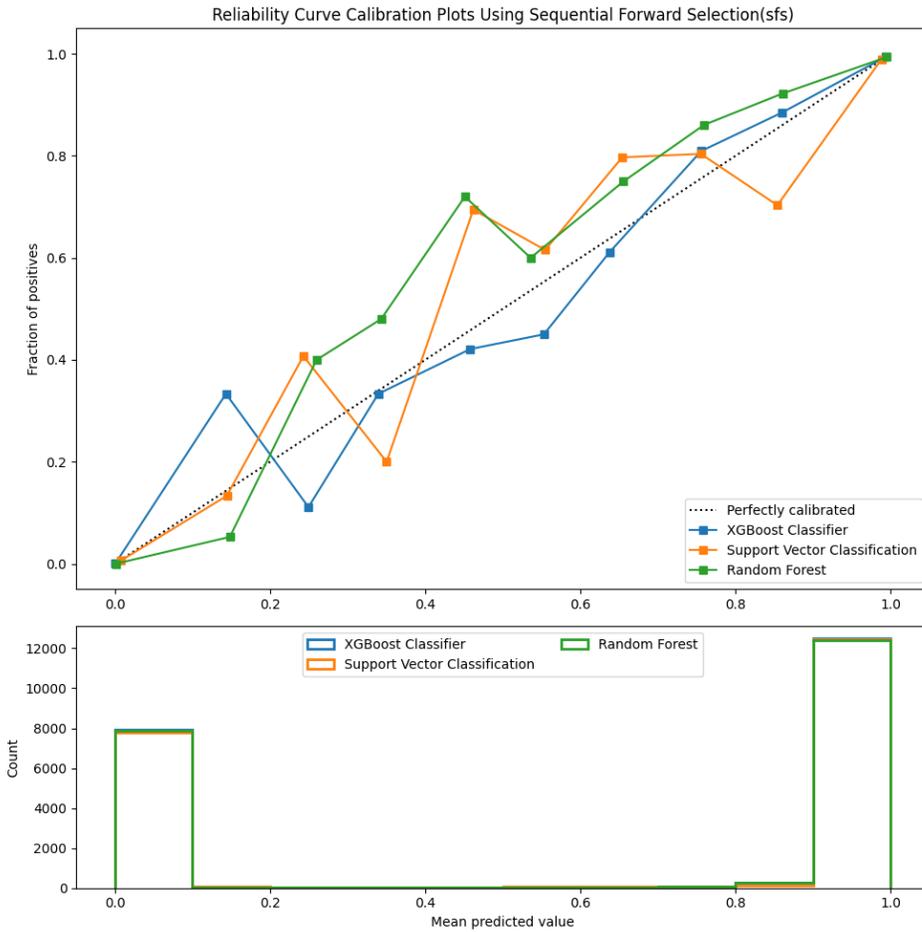**FIGURE 5:** ROC after applying sequential backward feature selection.

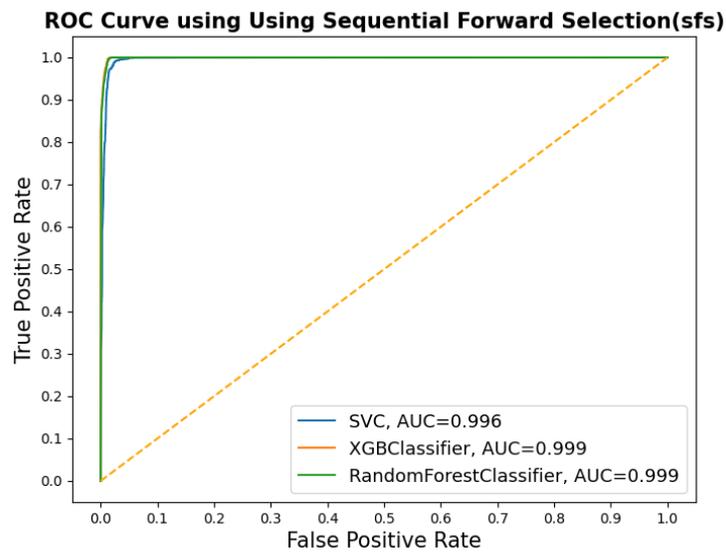**FIGURE 6:** Reliability Curve after applying sequential forward feature selection.



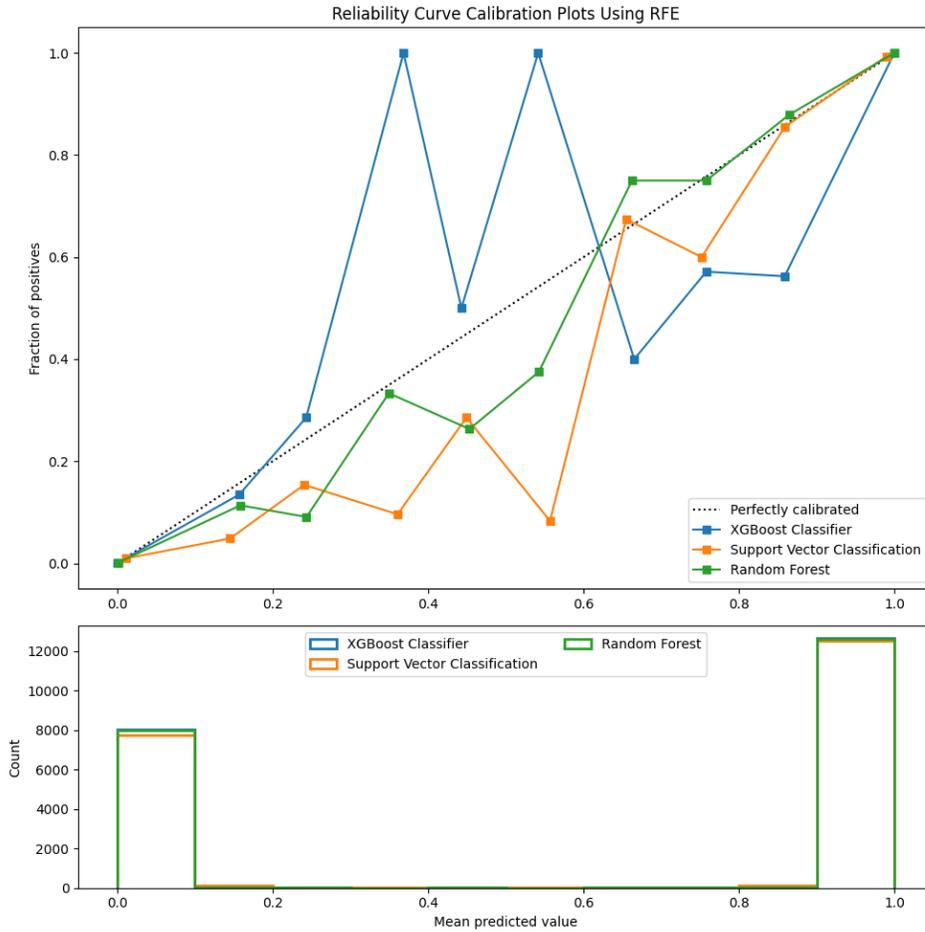**FIGURE 7:** ROC after applying sequential forward feature selection.

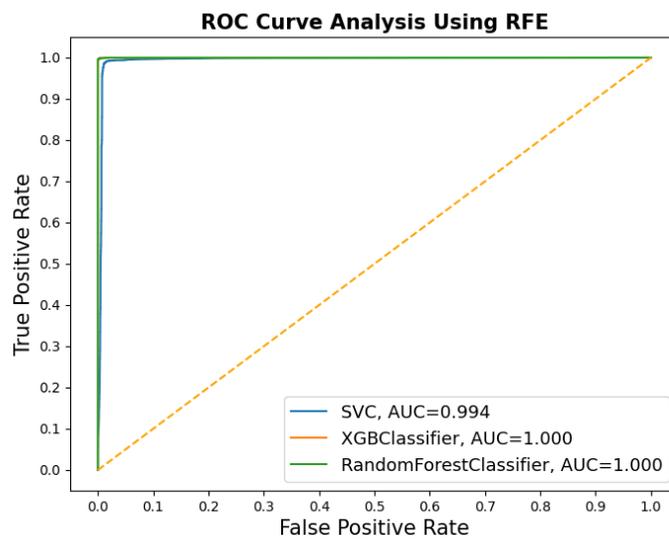**FIGURE 8:** Reliability Curve after applying RFE feature selection.



**FIGURE 9:** ROC after applying RFE feature selection.

## 7. CONCLUSION

Security threats to IoT systems and devices translate to significant security risks because of the inherent characteristics of the underlying technology. These characteristics make IoT environments versatile, functional, and efficient but can be vulnerable to threat actors. This research evaluates different supervised feature selection methods to predict malicious network traffic against IoT devices. We employ three different feature selection methods and implement three different logistic regression techniques for each selection method. We conclude that all the three logistic regression techniques (SVC, Random Forest, and XGBoost) performed with high accuracy. This implies that these techniques can be employed to predict an attack on IoT devices in a supervised learning setting.

Security attacks on IoT devices can sometimes be challenging to detect since IP addresses can be spoofed by the attacker, making it improper to be used as a machine learning feature. IP addresses are also mostly used in context with other indicators during intrusion detection. While this research can accurately predict an IoT attack, it should be noted that supervised learning is limited to the quality of training data and features selected. Statistical measures for feature selection must be carefully chosen and can significantly impact attack/intrusion predictions. The choice of limiting selection to eight features and limiting the dataset used for the month of May/2019 was purely for the research study. A different choice can result in different classifier accuracy results. Lastly, security against any asset should always be deployed in layers following risk, vulnerability, and threat analysis. A proactive effort by both manufacturers and the business community towards leveraging existing Cybersecurity controls, technology, and industry best practices frameworks can significantly mitigate the fast-rising IoT incident exposures.

## 8. ACKNOWLEDGEMENT

## 9. REFERENCES

[1] "New Cisco Annual Internet Report Forecasts 5G to Support More Than 10% of Global Mobile Connect — The Network." Internet: https://newsroom:cisco:com/pressrelease-content?type=webcontentfn&garticleId=2055169, [Jan. 05, 2021].

[2] S. Morgan, "Cybercrime To Cost The World $10.5 Trillion Annually By 2025," Internet: https://cybersecurityventures:com/hackerpocalypse-cybercrime-report-2016/fn#g:f_g:text=Cyberattacksarethefastestgrowing;insizefn%g2Csophisticationandcost:fn &gtext=oeDDoSattacksfn%g2Cransomwarefn%g2Cand;SharkonABC'sSharkTank, 2020, [Jan. 05, 2021].

[3] NIST Glossary, Internet: https://csrc.nist.gov/glossary/term/Cyber_Attack, [Jan. 05, 2021].

[4] "The "Only" Coke Machine on the Internet.", Internet: https://www:cs:cmu:edu/f_gcoke/history long:txt, [Jan. 05, 2021].

[5] Avast Threat Intelligence Team, "Hacker creates seven new variants of the Mirai botnet", Internet: https://blog:avast:com/hacker-creates-seven-new-variants-of-the-mirai-botnet, 2018, [Jan. 05, 2021].

[6] "How 5G and IoT devices open up the attack surface on enterprises - Security Boulevard.", Internet: https://securityboulevard:com/2020/04/how-5g-and-iot-devicesopen-up-the-attack-surface-on-enterprises/, 2020, [Jan. 05, 2021].

[7] P. A. Networks, "2020 Unit 42 IoT Threat Report", Internet: https://start:paloaltonetworks:com/unit-42-iot-threat-report, 2020, [Jan. 05, 2021].

Sundar Krishnan, Ashar Neyaz & Qingzhong Liu

[8]    I. Ullah and Q. H. Mahmoud, "A Scheme for Generating a Dataset for Anomalous Activity Detection in IoT Networks," in Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), vol. 12109 LNAI. Springer, may 2020, pp. 508–520. [On-line]. Available: https://link:springer:com/chapter/10:1007/978-3-030-47358-7 52, [Jan. 05, 2021].

[9]    H. Kang, D. H. Ahn, G. M. Lee, J. D. Yoo, K. H. Park, and H. K. Kim, "IoT network intrusion dataset", Internet: https://ieee-dataport:org/open-access/iot-network-intrusion-dataset, 2019, [Jan. 05, 2021].

[10]   "What is an Intrusion Detection System?", Internet: https://www:barracuda:com/glossary/intrusion-detection-system, [Jan. 05, 2021].

[11]   A. Kumar, W. Glisson, and H. Cho, "Network Attack Detection Using an Unsupervised Machine Learning Algorithm" in Proc. 53rd Hawaii Int.Conf. Syst. Sci. Hawaii International Conference on System Sciences, 2020, [On-line], Available: https://aisel.aisnet.org/hicss-53/st/cyber_threat_intelligence/8/, [Jan. 05, 2021].

[12]   S. Mukherjee, H. Asnani, E. Lin, and S. Kannan, Jul 2019, "ClusterGAN: Latent space clustering in generative adversarial networks," in 33rd AAAI Conf. Artif. Intell. AAAI 2019, 31st Innov. Appl. Artif. Intell. Conf. IAAI 2019 9th AAAI Symp. Educ. Adv. Artif. Intell. EAAI 2019, [On-line]. vol. 33, no. 01. AAAI Press, pp. 4610–4617.. Available: www:aaai:org, [Jan. 05, 2021].

[13]   L. F. Carvalho, S. Barbon, L. D. S. Mendes, and M. L. Proença, Jul 2016, "Unsupervised learning clustering and self-organized agents applied to help network management", Expert Syst. Appl., [On-line]. vol. 54, pp. 29–47, Available: https://dl.acm.org/doi/abs/10.1016/j.eswa.2016.01.032, [Jan. 05, 2021].

[14]   J. Sakhnini, H. Karimipour, and A. Dehghantanha, Aug 2019, "Smart Grid Cyber Attacks Detection Using Supervised Learning and Heuristic Feature Selection," in Proc. 2019 7th Int. Conf. Smart Energy Grid Eng. SEGE 2019. Institute of Electrical and Electronics Engineers Inc., [On-line] pp. 108–112., Available: https://ieeexplore.ieee.org/document/8859946, [Jan. 07, 2021].

[15]   L. Xiao, X. Wan, X. Lu, Y. Zhang, and D. Wu, Sep 2018, "IoT Security Techniques Based on Machine Learning: How Do IoT Devices Use AI to Enhance Security?" IEEE Signal Process. Mag., [On-line] vol. 35, no. 5, pp. 41–49, Available: https://ieeexplore.ieee.org/document/8454402, [Jan. 07, 2021].

[16]   M. Caron, P. Bojanowski, A. Joulin, and M. Douze, Jul 2018, "Deep Clustering for Unsupervised Learning of Visual Features," Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), [On-line] vol. 11218 LNCS, pp. 139–156, Available: http://arxiv:org/abs/1807:05520, [Jan. 07, 2021].

[17]   E. Balkanli, J. Alves, and A. N. Zincir-Heywood, Jan 2014, "Supervised learning to detect DDoS attacks," in IEEE SSCI 2014 2014 IEEE Symp. Ser.Comput. Intell. - CICS 2014 2014 IEEE Symp. Comput. Intell. Cyber Secur. Proc. Institute of Electrical and Electronics Engineers Inc., Available: https://ieeexplore.ieee.org/document/7013367, [Jan. 07, 2021].

[18]   V. Morfino and S. Rampone, Mar 2020, "Towards Near-Real-Time Intrusion Detection for IoT Devices using Supervised Learning and Apache Spark," Electronics, [On-line] vol. 9, no. 3, p. 444, Available: https://www:mdpi:com/2079-9292/9/3/444, [Jan. 09, 2021].

[19]   E. Hodo, X. Bellekens, A. Hamilton, P.-L. Dubouilh, E. Iorkyase, C. Tachtatzis, and R. Atkinson, 2016, "Threat analysis of IoT networks Using Artificial Neural Network Intrusion

Detection System" 016 International Symposium on Networks, Computers and Communications (ISNCC), [On-line] pp. 1-6, doi: 10.1109/ISNCC.2016.7746067. Available: https://ieeexplore:ieee:org/abstract/document/7746067/, [Jan. 09, 2021].

[20] C. Ioannou and V. Vassiliou, May 2019, "Classifying security attacks in IoT networks using supervised learning" in Proc. - 15th Annu. Int. Conf. Distrib. Comput. Sens. Syst. DCOSS 2019. Institute of Electrical and Electronics Engineers Inc., [On-line] pp. 652–658., Available: https://ieeexplore.ieee.org/abstract/document/8804727, [Jan. 11, 2021].

[21] E. Montalbano, "Report: Most Popular Home Routers Have 'Critical' Flaws," Internet: https://threatpost:com/report-mostpopular-home-routers-have-critical-flaws/157346/, 2020, [Jan. 11, 2021].

[22] S. Grimaldi, A. Mahmood, and M. Gidlund, Dec 2018, "Real-Time Interference Identification via Supervised Learning: Embedding Coexistence Awareness in IoT Devices," IEEE Access, [On-line] vol. 7, pp. 835–850, Available: https://ieeexplore.ieee.org/document/8570750, [Jan. 11, 2021].

[23] E. Anthi, L. Williams, M. Slowinska, G. Theodorakopoulos, and P. Burnap, Oct 2019, "A Supervised Intrusion Detection System for Smart Home IoT Devices," IEEE Internet Things J., [On-line] vol. 6, no. 5, pp. 9042–9053, Available: https://ieeexplore.ieee.org/document/8753563, [Jan. 11, 2021].

[24] I. El Naqa and M. J. Murphy, 2015, What Is Machine Learning? Cham: Springer International Publishing, [On-line] pp. 3–11, Available: https://doi:org/10:1007/978-3-319-18305-3 1, [Jan. 14, 2021].

[25] K. Kersting, Nov 2018, "Machine Learning and Artificial Intelligence: Two Fellow Travelers on the Quest for Intelligent Behavior in Machines," Front. Big Data, [On-line] vol. 1, p. 6, Available: https://www:frontiersin:org/article/10:3389/fdata:2018:00006/full, [Jan. 15, 2021].

[26] Y. Charfaoui, "Hands-on with Feature Selection Techniques: Wrapper Methods," Internet: https://heartbeat:fritz:ai/hands-onwith-feature-selection-techniques-wrapper-methods-5bb6d99b1274, 2020, [Jan. 15, 2021].

[27] "Recursive Feature Elimination — Yellowbrick v1.2.1 documentation.", Internet: https://www:scikit-yb:org/en/latest/api/modelselection/rfecv:html, [Jan. 15, 2021].

[28] "What is Supervised Learning?", Internet: https://www:ibm:com/cloud/learn/supervised-learning, [Jan. 16, 2021].

[29] "File:Random forest diagram complete.png - Wikimedia Commons.", Internet: https://commons:wikimedia:org/wiki/File:Randomforest diagram complete:png, [Jan. 16, 2021].

[30] "Caret: Classification and Regression Training.", Internet: https://cran:r-project:org/web/packages/caret/index:html, [Jan. 16, 2021].

[31] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, P. and Thirion, B. and Grisel, O. and Blondel, M. and Prettenhofer, A. and Weiss, R. and Dubourg, V. and Vanderplas, J. and Passos, and E. Cournapeau, D. and Brucher, M. and Perrot, M. and Duchesnay, "Scikit-learn: Machine Learning in fPgython." Internet: https://scikit-learn:org/stable/modules/generated/sklearn:svm:SVC:html, [Jan. 18, 2021].

[32] "Distributed Random Forest (DRF).", Internet: http://docs:h2o:ai/h2o/latest-stable/h2o-docs/data-science/drf:html, [Jan. 18, 2021].

[33]  T. Chen and C. Guestrin, Aug 2016, "XGBoost," Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discov. Data Min., Available: http://dx:doi:org/10:1145/2939672:2939785, [Jan. 18, 2021].

[34]  V. Morde and V. A. Setty, "XGBoost Algorithm: Long May She Reign!", Internet: https://towardsdatascience:com/https-medium-com-vishalmordexgboost-algorithm-long-she-may-rein-edd9f99be63d, 2019 [Jan. 18, 2021].

[35]  S.  Pafka,  "Benchmarking  Random  Forest  Implementations",  Internet: http://datascience:la/benchmarking-random-forestimplementations/, 2015, [Jan. 18, 2021].

[36]  K.  P.  Shung,  "Accuracy,  Precision,  Recall  or  F1?",  Internet: https://towardsdatascience:com/accuracy-precision-recall-orf1-331fb37c5cb9,  2015,  [Jan. 18, 2021].

Brindavana Sachidanand, Yuanyuan Jiang & Ahmad Hadaegh

# Twitter Based Sentimental Analysis of Impact of COVID-19 on Economy using Naïve Bayes Classifier

**Brindavana Sachidanand**                                   *sachi002@cougars.csusm.edu*
*Computer Science and Information System*
*California State University San Marcos*
*San Marcos, 92096, USA*

**Yuanyuan Jiang**                                           *yjiang@csusm.edu*
*Computer Science and Information System*
*California State University San Marcos*
*San Marcos, 92096, USA*

**Ahmad Hadaegh**                                            *ahadaegh@csusm.edu*
*Computer Science and Information System*
*California State University San Marcos*
*San Marcos, 92096, USA*

## Abstract

COVID-19 outbreak brought unprecedented changes to people's lives and made significant impact on the US and world economy. It wrought havoc on livelihood, businesses and ultimately the economy. Understanding how the sentiment on economy is changing and main factors that drives the change will help the public to make sense of the impact and generating relief measures. In this paper we present a novel Naïve Bayes model using a word-based training approach to perform the analysis and determine the sentiment of Twitter posts. The novelty of this methodology is that we use labelled set of words to classify the tweets to perform sentimental analysis as opposed to the more expensive methods of manually classifying the tweets. We then perform analysis on the resulting labelled tweets to observe the trend of economy from February 2020 to July 2020 and determine how COVID-19 impacted the economy based on what people posted on Twitter. We found our data was largely inclined towards negative sentiment indicating that the economy had been largely negatively impacted as a result of COVID-19. Further, we correlate the sentiment with the stock market index aka Dow Jones Industrial Average (DJIA) because stock market movement closely mirrors the economic sentiment and is shown as one of the main factors influencing people's attitude change from our sentimental analysis. We found strong correlation between the two, indicating stock market change is one of the driving factors behind people's opinion change about economy during pandemic. This work proposed and tested a generic lower-cost text-based model to analysis generic public's opinion about an event which can be adopted to analyze other topics.

**Keywords:** Sentimental Analysis, Naïve Bayes Classifier, COVID-19, Economy, Stock Market Index.

## 1. INTRODUCTION

COVID-19 (Corona Virus Disease '19) outbreak is unprecedented in human history and has disturbed the day-to-day life in different ways [1]. As a result of the pandemic, supply chain and people's economic behavior has deeply been affected. Manufacturing, transport, imports, and exports are greatly affected which posts impact to the economy [2]. People's lives and economic behavior are drastically different from pre-pandemic time as well. We are interested in studying how general public's opinion about the economy change during the pandemic and the main influential factors behind the change. By studying this topic, we can observe the extent to which

the economy has been impacted due to the pandemic from the general public's perspective, make meaningful observations and analysis which can help shine light on people's opinion about economic change and the main driving factors behind.

Though different techniques and data can be used to study the economy change, for example, unemployment rate, industrial production, etc. We are primarily interested in studying how it impacted personal lives. Twitter seems to be most reliable source to get the data [3], since Twitter reflects the sentiments of common people in day-to-day life and their views on various issues, whereas newspapers, news channels, magazines, etc. are the viewpoints of the journalists. Through Twitter we have access to analytics using which views can be statistically analyzed. Thus, collecting the data from Twitter is an effective and commonly used approach. So far, the existing works make use of already labelled tweets to perform training and get the sentiments [4]. Labelled tweets refer to the classified tweets i.e., tweets along with their sentiment. Manually labelling large amount of tweets can be very resource intensive.

In this paper, we present a novel methodology to train data based on the labelled set of words i.e., word-based training to classify the tweets to perform sentimental analysis instead of the more expensive method of manually classifying the tweets. Once this model is built, it can be used on any text data to perform sentimental analysis. In this approach, we use Twitter data as an example to showcase our model and study the impact of COVID-19 through sentimental analysis. This paper focuses on public sentiment by collecting the tweets having keywords like 'economy', 'stock market' and 'GDP' since the pandemic started.

In section 2, we outline the background work and also discuss the related work in detail. We have collected tweets from February 2020 to July 2020 and trained our model using Naïve Bayes classifier and further used the classifier to perform Sentimental Analysis [5]. The classifier training is based on the set of positive words and negative words [6]. Once the model is trained, we use the trained classifier to calculate positive and negative probability of a tweet and accordingly classify the tweet as positive, negative, or neutral. This is explained in detail under section 3. After determining sentiments for every tweet, we observe and analyze the obtained results to understand the impact of COVID-19 on economy and people's opinion change. From the results, we found stock market change seems to be one of the main factors driving people's opinion change about the economy. We also study the correlation of the inferred sentiment with the stock market index i.e., Five-day moving Dow Jones Industrial Average (FDJIA) [7] because it is generally believed that the stock market is a reflection of the current state of economy and the movement of stock market captures the movement of economic sentiment. This is explained in detail under section 4. Lastly, we conclude our work in section 5.

## 2. BACKGROUND AND RELATED WORKS

Social Media has become an extension of who we are, all that we post on social media is related to our feelings and our opinions about a particular matter in hand [8]. It is a place where large amount of data gets generated continuously. Twitter is one of the popular social media platforms where people are open to share their thoughts and concerns about any ongoing affairs, matters or crisis [1]. It is a rich source of information that would make a path for the analysis of social phenomena and its related sentiment [3]. Twitter data is the best source to do sentimental analysis (mainly categorized as positive, negative and neutral) on COVID-19's impact on Economy. Analyzing the tweets gives us insights into user's expressions, opinions and attitude [5].

Traditional approach to estimate sentiments involves training on already labelled data. This training is further used on the actual tweets to infer the sentiments. Whereas, in this work we

Brindavana Sachidanand, Yuanyuan Jiang & Ahmad Hadaegh

propose a novel approach of word-based training where we use certain key words to determine the sentiments of the tweets. We use a set of positive words and negative words, combine them, and use them as training data for the Naïve Bayes model. Once the model is trained, we obtain a trained classifier to get the sentiment of a tweet. We do not use already labelled tweets which is a very labor intensive approach. Instead, we make use of labelled words to label the tweets using Naïve Bayes approach. This is a novel approach which is intuitively modular and yields good accuracy as demonstrated by the testing results.

There is some earlier work done on Twitter sentimental analysis using Naïve Bayes approach. In one approach, tweets that are already labelled as positive or negative are taken as the training set [9]. Number of positive words and negative words are counted for each set. If a word exists in only the positive set, the word has positive weight. If a word exists in only the negative set, the word has negative weight. Otherwise, if a word exists in both sets, it has both positive and negative weight [9]. In another method, the weights are modified using the average of weight differences [9].

COVID-19 impact on economy has been discussed in other works. In A. Atkeson 's work [10], a SIR Markov model is built wherein the population is divided into 3 categories: susceptible to disease (S), actively infected with disease (I) and recovered (R) [10]. Importance is given when the fraction of active infections in the population exceeds 10% which can result in poor economic condition and cumulative impact of the disease over an 18-month horizon. Through this, they draw conclusions of the impact of COVID-19 on economy to be adverse and suggest that mitigation measures are needed.

B. Le and H. Nguyen used SVM and Naïve Bayes classifiers to categorize tweets into positive or negative [11]. They achieved an accuracy of 80.00% using Naïve Bayes approach and 78.08% using SVM approach. Bishwo Prakash Pokharel used TextBlob library of Python to classify tweets into positive and negative [12].

A Agarwal et al [13] Performed Sentimental Analysis of Twitter data using the POS-specific prior polarity features and a tree kernel. They demonstrated that tree kernel and feature-based approach outperform the unigram baseline using SVM. They achieved an accuracy of 75.39%.

Classifiers like Naïve Bayes, Maximum Entropy, and Support Vector Machines were used to classify the tweets as positive and negative by A. Go, L. Huang and R. Bhayani [14]. They use tweets with emoticons for distant supervised learning. The accuracy achieved in the case of Naïve Bayes classification was 81.30%.

COVID-19 has sparked interest in analysis of the impact through Twitter and sentiment analysis. KH Manguri et al measured how COVID-19 was trending on Twitter [15], Sakun Boon-Itt et al discovered the top three concerns related to COVID-19 on Twitter through sentiment analysis [16]. Kalifer Garcia et al did sentiment analysis on COVID-19 related tweets and concluded that the sentiment was largely negative [17]. Usman Naseem et al did Twitter based sentiment analysis on COVID-19 related tweets in early stages of pandemic and concluded that the sentiments were negative [18]. V. Senthil et al. used Twitter sentiment analysis to study the impact of COVID-19 on travel industry [19].

Table 1 shows how our work is different from the already existing work based on two features - training and impact of COVID-19 on economy. This table shows the novelty of our work.

| Feature | Existing Work | Our Work |
|---|---|---|
| Training | In most cases, the tweets are already labelled. Training and testing are done on the model based on these training and testing sets [4]. | Instead of using the manually labelled tweets, we rather use labelled words (words labelled as positive or negative), and the training and testing are based on set of positive words and negative words (word-based training). We then predict the sentiment of a tweet based on the resultant trained classifier. |
| Showcasing impact of COVID-19 on economy | Many theoretical surveys were done on the impact of COVID-19 on economy [10]. There are also few models like SIR Markov Model that are used to relate the number of increasing active COVID-19 cases to economic impact [11]. | No direct analysis of economy using computational methods exists in existing work so far. We have used Naïve Bayes Classifier approach on tweets to do sentimental analysis by collecting tweets related to COVID-19 and economy and further corelating the sentiments with stock market index |

**TABLE 1:** Differences between existing work and our work.

## 3. METHODOLOGY

The entire implementation is done using Python 3., we collect the data from Twitter which is one of the most popular social media platform where people express their views. There are other social media platforms like Facebook, LinkedIn, Reddit exist, and it might be possible to collect data and perform similar sentimental analysis through other platforms if the data is text based and can be cleaned into similar format. But Twitter is the most suitable platform for our purpose since tweets are mostly text based with large number of active users expressing personal views in real time. We train our model using Naïve Bayes classifier and further use the classifier to perform Sentimental Analysis [6]. The Python libraries we use are nltk's NaiveBayesClassifier (to train the model), re (to clean the tweets), matplotlib (to plot graphs), WordCloud (to draw a word cloud having most frequently used words), pandas (to create a data frame), LangDetect (to fetch only English tweets). We start by collecting tweets related to economy, clean the tweets, and tokenize them. Rather than manually labelling the tweets as positive or negative, we use a novel approach using word level labelling. In this approach, we have a set of positive words and negative words stored in separate files. We store the list of words in a list where each word is stored along with the sentiment as positive or negative. The tweets whose sentiment need to be predicted are tokenized since the trained data set (the set of positive and negative words) is in tokenized form. Then, we train our model using Naïve Bayes classifier on the word set that was obtained from the set of positive words and negative words. Once the model is trained, we use the trained classifier to calculate positive and negative probability of a tweet and accordingly classify the tweet as positive or negative. We also have neutral tweets, that will be explained in detail under Step 9 of section 3.4. After determining sentiment for every tweet, we observe the results, and do the analyses on how the tweets' sentiments vary from month to month and correlate the percentages of positive and negative tweets with Dow Jones Industrial Average (DJIA) to observe how the stock market index correlates to the sentiment of tweets. The details are explained in section 4. Figure 1 shows the steps involved in our approach.
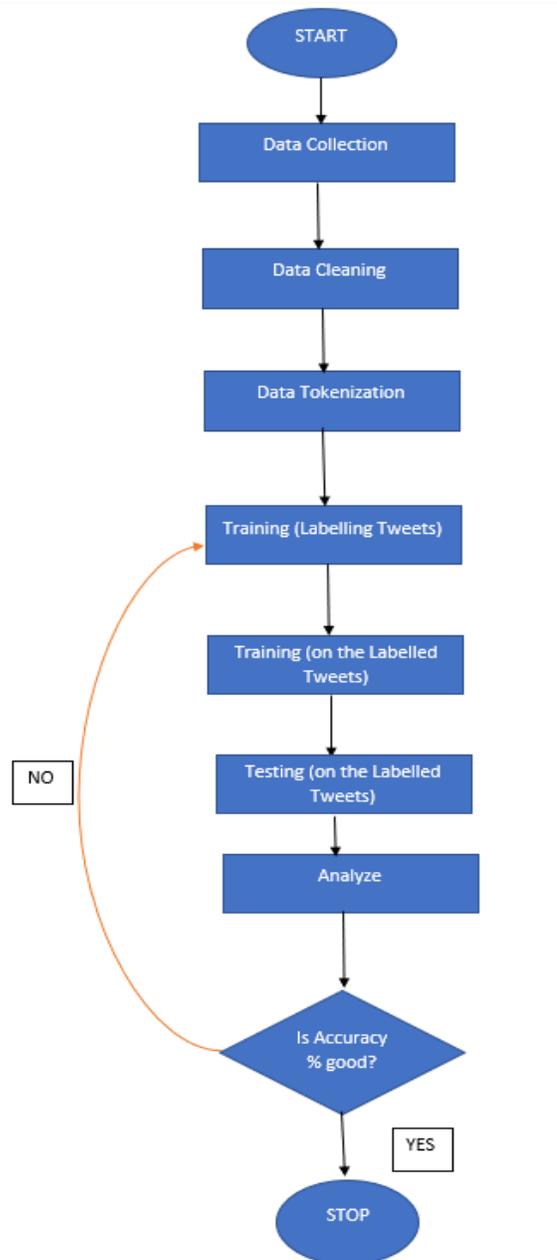
Brindavana Sachidanand, Yuanyuan Jiang & Ahmad Hadaegh



**FIGURE 1:** Steps involved in the approach.

### 3.1 Data Collection
Initially, we started collecting tweets using "Tweepy" streaming API. But since February 1, 2018, we cannot access tweets that are older than 7 days. So, we used GOT3 API (Get Old Tweets). GOT3 is a GitHub repository which uses "URLlib" for fetching tweets from Twitter's advance search. We used the search queries as *'economy', 'stock market', 'GDP', 'unemployment' and 'business shutting'* to collect the tweets. We collected about 2000 tweets per day and tweets range from February 2020 to July 2020. We focused on collecting tweets for 8 days in every month. We randomly selected some days from every month i.e., mostly weekends because users tweet the most during weekends in their down time. In total, we collected tweets for 48 days i.e., 6 months multiplied by 8 days, so approximately there were around 100,000 total tweets collected.

Each day's tweets are stored in a separate folder. We fetched only English tweets using "langdetect" library from Python and setting language to 'en'. Example 1 shows an instance of a collected raw tweet.

*Example 1: "@SenatorWong so good to see that the Libs have an extra $60 billion to put back into the economy. Rather than be at a loss? The Libs have managed COVID exceptionally well. Labor are trying to find fault. Won't work. Labour lack finesse and ability to even build wealth."*

### 3.2 Data Cleaning
Once data is collected, we created data frame for the tweets using "Pandas" library of Python. Next step in the process is to clean the tweets. Example 2 shows an instance of a cleaned tweet. Tweets are cleaned using regular expressions in Python that are available through the "re" library. The usernames, hashtags, and links were removed as these were unnecessary tokens while analyzing a tweet. Usernames and links are unnecessary to determine the sentiment of a tweet and hence we remove them. Hashtags helps expressing the meaning or theme of a tweet. But they may not always be useful for determining the sentiment of a tweet since not every tweet uses them and the Hashtags might not cover the entire meaning of the message. Thus, they are removed during data cleaning.

*Example 2: "so good to see that the Libs have an extra $60 billion to put back into the economy Rather than be at a loss The Libs have managed COVID exceptionally well Labor are trying to find fault Won't work Labour lack finesse and ability to even build wealth."*

### 3.3 Data Tokenization
Cleaned tweets are passed to a tokenize function to remove stop words, punctuations and finally tokenize every tweet. The result of a tokenized tweets is the list of all the tokens appearing in each tweet separated by commas. Example 3 shows an instance of a tokenized tweet. It is necessary to get each word in a tweet as a token because in the training step we would be training on a set of positive and negative words (or tokens), and then apply the tokenized tweets on the trained classifier.

*Example 3: "['good', 'see', 'libs', 'extra', '60', 'billion', 'put', 'back', 'economy', 'rather', 'loss', 'libs', 'managed', 'covid', 'exceptionally', 'well', 'labour', 'trying', 'find', 'fault', 'work', 'labour', 'lack', 'finesse', 'ability', 'even', 'build', 'wealth']"*

### 3.4 Training (For labelling the tweets)
We construct 2 files containing the positive and negative words. One is positive words file which has a list of various positive words and another is negative words file which has a list of various negative words. The key to classify our tweets as positive or negative is to perform the training using Naïve Bayes classifier on the labelled set of words and once trained, use the trained classifier to find positive and negative probability of every tweet. Then based on the results, we classify the tweets as either positive, negative or neutral [5]. We use nltk's NaiveBayesClassifier library in Python in this step [20]. This includes 10 steps as discussed below.

**Step 1:** We constructed 2 word files – one having a list of positive words (about 2000 words) and another having a list of negative words (about 5000 words) and converted the words into a word dictionary.

**Step 2:** Once the word dictionary is created, we combine both positive and negative word dictionaries and store it in a variable named training. Example 4 shows a snapshot of the tokenized dictionary generated after combining both the dictionaries. This word dictionary constructed from the word files would further help us in classifying a tweet as positive or negative while training. Notice that each element is a list containing a dictionary and a string index i.e., either positive or negative.

*Example 4: [{'abound': True}, 'positive'], [{'abounds': True}, 'positive'],
[{'abundance': True}, 'positive'], [{'abundant': True}, 'positive'], ……..,
[{'wrong': True}, 'negative'], [{'wrongful': True}, 'negative']*

**Step 3:** Naïve Bayes Classifier is used to train the model on the training data
[8].

**Step 4:** The Naïve Bayes training generates a classifier. This classifier will
be used in the further steps to decide whether the tweet is positive or
negative [5].

**Step 5**: Now we go back to our tweet. Recall that we tokenized our tweet.
We apply a function on the tokenized tweet and generate a word dictionary
of the tokenized tweet in the same way that we performed for the set of
positive and negative words. Example 5 shows an instance of a tokenized
tweet dictionary. Notice that each token is a key-value pair.

*Example 5: "{'good': True, 'see': True, 'libs': True, 'extra': True, '60': True,
'billion': True, 'put': True, 'back': True, 'economy': True, 'rather': True, 'loss':
True, 'libs': True, 'managed': True, 'covid': True, 'exceptionally': True, 'well':
True, 'labour': True, 'trying': True, 'find': True, 'fault': True, 'work': True,
'labour': True, 'lack': True, 'finesse': True, 'ability': True, 'even': True, 'build':
True, 'wealth': True}"*

**Step 6:** Since the trained classifier is trained using Naïve Bayes, it is based
on probabilities [9]. Now a positive and negative probability score would be
assigned to the tokenized tweet dictionary that was generated in Step 5. The
probability value ranges from 0 to 1. Few operations are applied on the
classifier that is  generated by the Naïve Bayes Classifier to get the positive
and negative probability score of the tokenized tweet dictionary.

- prob_result = classifier.prob_classify (tokenized_tweet_dictionary)
- positive_prob = prob_result.prob("positive")
- negative_prob = prob_result.prob("negative")

**Step 7:** positive_prob and negative_prob give us the positive and negative
probability score of the tokenized tweet dictionary.

**Step 8:** Example 6 shows how a sentiment of a tweet is decided. If the
positive probability is greater than the negative probability, then the tweet is
classified to be positive, else it is a negative tweet.

**Step 9:** According to our algorithm, we have classified all the tweets that
have positive and negative probability ranging between 0.4 and 0.6 to be
neutral tweets. Thus, we end up with positive, negative, and neutral tweets.

*Example 6:*
- *Positive Probability: 0.9060365743100163,*
- *Negative Probability: 0.09396342568998274*
- *Sentiment: Positive*

In Example 6, the Positive Probability of the tweet is greater than the
Negative Probability. Hence, the tweet is Positive. If the probabilities ranged
between 0.4 and 0.6 the tweet would have been Neutral.

**Step 10:** Since we obtained the sentiment of every tweet, percentages of positive, negative, and neutral tweets are calculated for each day.

We have collected tweets for 8 days in every month from February to July. For each day, we get cleaned tweets, tokenized tweets, sentiment for every tweet and percentages of positive, negative, and neutral tweets.

After training whenever we observed that a particular tweet is not classified correctly, we went back to the positive and negative word files and checked if the required word is present in the word files. In case of a word not being present, we manually added the required word to the respective file (positive or negative file depending on the case). We kept adding new words to the files as and when required and kept re-training and retuning the model. This enabled us to achieve better training and better classification of the tweets.

We collected around 100,000 tweets in total to perform sentimental analysis. We feel that this number of tweets is good to perform the analysis. However, a greater number of tweets could be collected to perform the analysis as well, there is no inherent limitation in the methodology.

### 3.5 Training & Testing on the labelled tweets
After generating the labelled tweets, we split the labelled tweets into training and testing data i.e., 80% of the labelled tweets would be training data and 20% of the labelled tweets would be testing data [8]. We pass these training data and testing data to the Naïve Bayes Classifier to calculate the accuracy. We have achieved an accuracy of 85.92 %.

The key to achieve good training is to include as many positive and negative words as possible in the training set that we used for labelling the tweets so that they cover almost all words present in the tweets and accordingly classify the tweet as either positive or negative based on those words. If this initial labelling of tweets happens well, then our model constructed is effective. We performed the word-based training with some corrections multiple times in order get good labelling of tweets. Thus, good labelling of tweets would ensure good accuracy measure when the labelled tweets are tested.

### 3.6 Analysis
The sentimental analysis results obtained from the model is used to draw bar plots separately for every month to observe how the sentiments differ on particular days in every month. Plots are drawn using Python's matplotlib library. Figure 2 shows the plot having sentimental analysis from Feb 2nd, 2020 to July 19th, 2020. X axis gives the days and Y axis gives the percentages of positive, negative, and neutral tweets.

From Figure 2 it can be observed that the tweets are mostly inclined towards positive in February and more towards negative from March through July. Thus, the data found indicates that people have more negative tweets than positive tweets related to economy. We are interested in studying the underlining reasons that cause this effect. Since the collected tweets are related to economy, our hypothesis is that "bad economy might be one of the top reasons that made people generate negative posts". We will support our hypothesis with explanation in section 4.

Python's wordcloud library is used to draw a word cloud from the tokens obtained through the tweets to observe the most important and frequently used words which would depict the relevant words related to the current ongoing situation. Those words that appear in larger font in the word cloud depict that they are the mostly occurring words and the ones in the smaller font are the least occurring words. For e.g., the word 'stock market' is the most frequently used word according to the word plot. Thus, stock market is most frequently discussed topic during February 2020 to July 2020 which might indicate that Covid-19, among other factors, brought changes and people's attention to the stock market. We can also observe that there are many more frequently

occurring words. Figure 3 shows the Word Cloud obtained from the tweets.



**FIGURE 2:** Plot showing sentimental analysis from February to July 2020



**FIGURE 3:** Word Cloud.

## 4. RESULTS

We have presented a novel approach for sentimental using word-based training based on Naïve Bayes approach. As explained in section 3.5, this approach yields good accuracy (85.92 %) in testing. Table 2 shows comparisons of accuracies across related works and our work.

| Works | B. Le Work [11] | A. Agarwal Work [13] | Alec Go Work [14] | Our Work |
|---|---|---|---|---|
| **Accuracy** | 80.00% | 75.39% | 81.30% | 85.92% |

**TABLE 2:** Comparisons of accuracies across related works.

As seen in Table 2, our approach gives accuracy which is in the higher range of the results. Thus, these results validate our model. In sections 5.1 and 5.2 we further analyze the results obtained from our model.

**4.1 Results from Monthly Sentimental Analysis**
From the plots obtained, it can be observed that our data is mostly inclined towards negative sentiment while being positive on few days. When we observe the plots obtained in each month specifically, we observe that in the month of February, the sentiment is mostly positive and rarely negative, from March through May the sentiment is highly negative and from June through July the positive sentiment percentage gradually increases while still being negative on few days. We can validate the trend observed by quoting that in February the sentiment is mostly positive because the COVID-19 pandemic had not yet set in and the world economy was in a good shape. But during March, the pandemic set in and the economic conditions started to worsen. Hence, the sentiments become more negative and less positive. From March through May the sentiment is mostly negative as the economic conditions became worse because of lockdown and stay at home orders. From June through July the government started taking necessary steps to improve the economy and in addition, many states started opening up gradually towards the end of June. As the result, the sentiment started inclining towards positive while still retaining negative sentiments on few days.

Thus, our earlier hypothesis stated under section 3.6 that "bad economy might be one of the top reasons that made people generate negative posts" is validated. Since we collected tweets using the keywords *'economy', 'stock market', 'GDP', 'unemployment' and 'business shutting'* the sentiment calculated can be directly attributed to the state of economy at that point of time.

It is observed that positive and negative tweets have higher sentiments (around 40% to 50%) and neutral percentage is usually lower (around 10%). Positive and negative percentages have around 10% difference because whenever the economy drops, the government comes up with remedial measures to make the positive percentage to be around 10% closer to the negative percentage. The tweets are mostly either positive or negative.

**4.2. Correlating percentages of positive and negative tweets with Dow Jones Industrial Average (DJIA)**
Now that we have sentimental analysis in place, we can move on to correlating the sentiments to stock market. It is generally believed that the stock market is a reflection of the current state of economy and the movement of stock market captures the economic sentiment. Dow Jones Industrial Average (DJIA) is a USA stock market index that consists of 30 blue chip publicly traded companies. This largely tracks the US stock market sentiment [7]. Five-day moving Dow Jones Industrial Average (FDJIA) gives an average of Dow Jones Industrial Average over 5 days [7].

All tweets collected are in English and since almost 85% of the tweets are from United States, we would focus on Dow Jones Industrial Average specifically in United States to draw correlations to percentage changes in positive and negative tweets. Correlation is a statistical measure which indicates how two or more entities are related to each other [21]. A positive correlation indicates that the entities are directly related and shows the measure of how the entities increase or decrease in parallel; a negative correlation indicates that the entities are inversely related and shows the measure of how one entity increases as the other decreases [22].

In our case, it is observed that the economic sentimental analysis is corelating well to DJIA [7]. Also, since the word cloud seems to suggest that there is a good correlation between the stock and sentiments, we tried to correlate the stock index with the sentiment – it seems to corelate reasonably well.

We calculate the Five-day moving Dow Jones Industrial Average (FDJIA). It is seen that whenever the percentage change in FDJIA increases, there is an increase in change of positive percentage of tweets and decrease in change of negative percentage of tweets. Basically, this means that whenever the stock market is going down, the sentiment of tweets is mostly negative i.e., the negative sentiment percentage increases. Similarly, whenever the stock market is going

up, the sentiment of tweets is mostly positive i.e., the positive sentiment percentage increases [21]. Figure 4 shows the Relationship between percentage change in FDJIA with percentage change in positive tweets and percentage change in negative tweets . X axis gives the days and Y axis gives the percentages of positive, negative, and neutral tweets.

From Figure 4, it can be observed that whenever there is a rise in percentage change in FDJIA, there is a rise in percentage change in positive tweets i.e., they are positively (directly) corelated. But whenever there is a drop-in percentage change in FDJIA, there is a rise in percentage change in negative tweets i.e., they are negatively (inversely) correlated.
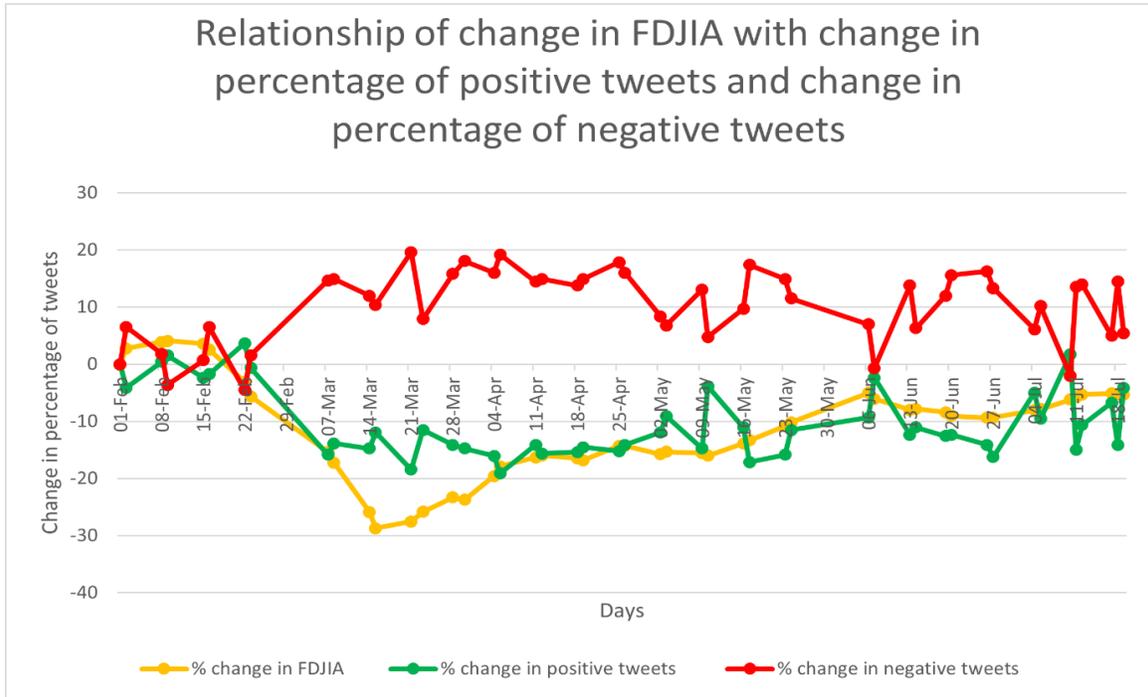


**FIGURE 4:** Relationship between percentage change in FDJIA with percentage change in positive tweets and percentage change in negative tweets.

Figure 5 (a) and (b) are the plots showing correlation between the respective entities. X axis gives the %change in FDJIA and Y axis gives the %change in positive and negative tweets in 5 (a) and 5 (b) respectively.
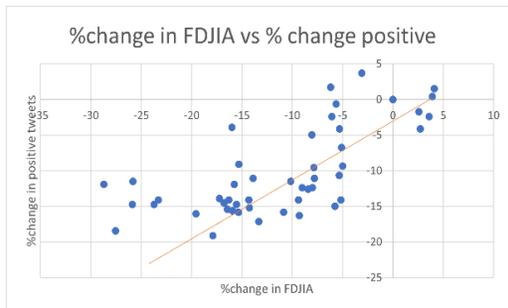


**FIGURE 5 (a):** Correlation between percentage change in FDJIA and percentage change in positive tweets.
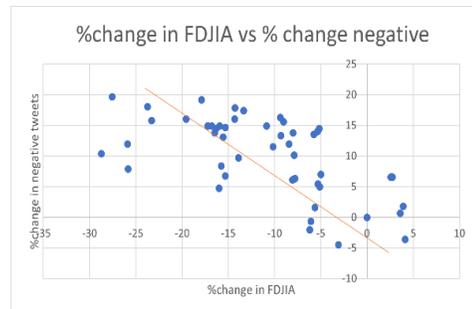


**FIGURE 5(b):** Correlation between percentage change in FDJIA and percentage change in negative tweets.

From Figure 5 (a) it can be observed that there is direct correlation (positive correlation) between percentage change in FDJIA and percentage change in positive tweets and similarly, from Figure 5 (b) it can be observed that there is inverse correlation (negative correlation) between percentage change in FDJIA and percentage change in negative tweets.

We use Equation 1 to calculate the correlation value between the entities [21].

$$r = \frac{1}{n-1} \sum \left(\frac{x - \overline{x}}{S_x}\right)\left(\frac{y - \overline{y}}{S_y}\right)$$

**Equation 1:** Correlation.

In Equation 1, n refers to the number of entries, x and y refer to the entities whose correlation need to be determined, $\overline{x}$ and $\overline{y}$ refer to mean/average of x and y, respectively. $S_x$ and $S_y$ refer to the Standard Deviations of x and y, respectively [22].

After calculating, we obtain correlation value as **0.69** in the case of 5 (a) and **-0.61** in the case of 5 (b). This proves that 5 (a) is the case of positive correlation and 5 (b) is the case of negative correlation.

Thus, it is observed that the percentage change in Five-day moving Dow Jones Moving Average (FDJIA) is positively correlated to percentage change in positive tweets and negatively correlated to percentage change in negative tweets. This means that whenever the stock market is performing well, there is an increase in the FDJIA, and this in turn reflects in the sentiments of tweets and thus people tweet positively. Similarly, whenever the stock market is not performing well, there is a decrease in the FDJIA, and this in turn reflects in the sentiments of tweets and thus people tweet negatively.

## 5. CONCLUSIONS

In this paper we presented a novel approach to do sentimental analysis using Naïve Bayes Classifier which is based on word-based training. We used this analysis on Twitter data to track the sentiments of the people on the effect of COVID-19 on world economy. We collected tweets related to economy during a 6-month period i.e., February 2020 to July 2020. We then cleaned, tokenized the tweets and then classified each tweet as positive, negative, or neutral using Naïve Bayes classifier. This novel approach is modular and yields good results. Once we have the labelled tweets i.e., each tweet along with its sentiment, we performed training (80% of labelled tweets) and testing (20% of labelled tweets) and found accuracy to be 85.92 %.

We noticed the impact of COVID-19 on economy from the obtained results and found that the tweets were inclined towards positive sentiment in February 2020, from March through May it mostly inclined towards negative sentiment and from June through July we saw a slight increase in positive sentiment while being negative on few days. Thus, we reached a conclusion that our data was largely inclined towards negative sentiment indicating that the economy had been largely negatively impacted as a result of COVID-19.

 We also showed strong correlation of the sentiment with Five-day moving Dow Jones Industrial Average (FDJIA) i.e., we observed positive correlation between percentage change in positive tweets and percentage change in FDJIA and negative correlation between percentage change in negative tweets and percentage change in FDJIA. This means that whenever there is positive sentiment in tweets, the FDJIA is also in positive territory indicating that the stock market is performing well. Similarly, whenever there is negative sentiment in tweets, the FDJIA is also in negative territory indicating that the stock market is not performing well.

This work proposed and tested a generic lower-cost text-based model to analysis generic public's opinion about an event which can be adopted to analyze other topics. Our model could be used to analyze sentiment on any piece of text be it any email, news article or a message. For example, a higher-level application can be built to classify any news article to help the reader with the sentiment of the article. Our model can also be used by the administration to gauge the sentiment of people on various issues that affects the people. It could be a good tool to gather sentiments of people if any policy change is done. It could also be used in election predictions where psephologists can use the tweets to predict the outcome of an election. For example, based on the sentiment of the tweets, one can get the approval ratings of candidates in an election - this can be useful tool for the candidates to decide on the strategy to be used in the course of their election campaign.

A limitation of our work could be that since not everyone uses Twitter, the sentiments of those cluster of people not using Twitter are not captured. We do not know if their sentiments conform to the results that we provided, or if they deviate. More experiments can be done using our methods with more tweets or tweets from different time to better analyze the accuracy and reliability of this method. It would be interesting to see if similar results can be found via posts from other platforms, like Facebook or Reddit.

The scope of our work was between the 6-month period from February 2020 to July 2020. As a future scope, we could also collect tweets for a period of 1 year i.e., from 2020 to 2021 and observe the overall 1-year trend. We can also factor in the quantity of likes, retweets, and shares in the form of weights which would be added to the raw sentiment calculated by the sentiment analysis. In our work we studied the relationship between the sentiment towards economy and stock market index - as a future work we could also consider relationship of sentiment towards economy with other factors such as number of COVID-19 cases, number of COVID-19 deaths and number of remedial policies introduced by the government.

## 6. REFERENCES

[1] S.Hamidian and M. Diab, "Rumor Detection and Classification for Twitter Data". SOTICS 2015 : The Fifth International Conference on Social Media Technologies, Communication, and Informatics.

[2] Cinelli, M., Quattrociocchi, W., Galeazzi, A. *et al.* The COVID-19 social media infodemic. *Sci Rep* 10, 16598 (2020). https://doi.org/10.1038/s41598-020-73510-5.

[3] Amedie, Jacob, "The Impact of Social Media on Society" (2015). *Pop Culture Intersections*. 2.

[4] Fang, X., Zhan, J. Sentiment analysis using product review data. Journal of Big Data 2, 5 (2015). https://doi.org/10.1186/s40537-015-0015-2

[5] Rish, Irina. (2001). An Empirical Study of the Naïve Bayes Classifier. IJCAI 2001 Work Empir Methods Artif Intell. 3.

[6] I.V., Shravan. (2016). Sentiment Analysis in Python using NLTK. OSFY - OpensourceForYou.

[7] R. Glen Donaldson and Harold Y. Kim, "Price Barriers in the Dow Jones Industrial Average",The *Journal of Financial and Quantitative Analysis*, Vol. 28, No. 3 (Sep 1993), pp. 313-330

[8]   C. Troussas, M. Virvou, K. J. Espinosa, K. Llaguno and J. Caro, "Sentiment analysis of Facebook statuses using Naive Bayes classifier for language learning," *IISA 2013*, Piraeus, 2013, pp. 1-6, doi: 10.1109/IISA.2013.6623713.

[9]   J. Song, K. T. Kim, B. Lee, S. Kim and H. Y. Youn, "A novel classification approach based on Naïve Bayes for Twitter sentiment analysis," KSII Transactions on Internet and Information Systems, vol. 11, no. 6, pp. 2996-3011, 2017. DOI: 10.3837/tiis.2017.06.011.

[10]  A. Atkeson, "What will be the Economic Impact of Covid-19 in the US? Rough Estimates of disease scenarios", NATIONAL BUREAU OF ECONOMIC RESEARCH 1050 Massachusetts Avenue Cambridge, MA 02138 March 2020

[11]  B. Le and H. Nguyen, "Twitter Sentiment Analysis Using Machine Learning Techniques," Advanced Computational Methods for Knowledge Engineering, pp.279-289, 2015.

[12]  B. P. Pokharel, "Twitter Sentiment analysis during COVID-19 Outbreak in Nepal", M.Phil in ICT 3rd semester Nepal Open University, Nepal, May 2020.

[13]  A. Agarwal, B. X, I. Vovsha, Owen Rambow and Rebecca Passonneau, "Sentiment analysis of Twitter data," in Proc. of the Workshop on Languages in Social Media, pp.30-38, 2011.

[14]  A. Go, L. Huang and R. Bhayani, "Twitter Sentiment Classification using Distant Supervision". 2009.

[15]  K. H. Manguri, R. N. Ramadhan, and P. R. Mohammed Amin, "Twitter Sentiment Analysis on Worldwide COVID-19 Outbreaks", Kurdistan Journal of Applied Research, vol. 5, no. 3, pp. 54-65, May 2020.

[16]  Boon-Itt S, Skunkan Y, "Public Perception of the COVID-19 Pandemic on Twitter: Sentiment Analysis and Topic Modeling Study", JIMR Public Health Surveill 2020;6(4):e21978.

[17]  Klaifer Garcia, Lilian Berton, "Topic detection and sentiment analysis in Twitter content related to COVID-19 from Brazil and the USA", Applied Soft Computing, Volume 101, 2021, 107057, ISSN 1568-4946

[18]  U. Naseem, I. Razzak, M. Khushi, P. W. Eklund and J. Kim, "COVIDSenti: A Large-Scale Benchmark Twitter Data Set for COVID-19 Sentiment Analysis," in *IEEE Transactions on Computational Social Systems*, doi: 10.1109/TCSS.2021.3051189.

[19]  Senthil V., Goswami S. (2020) An Exploratory Study of Twitter Sentiment Analysis During COVID-19: #TravelTomorrow and #UNWTO. In: Sharma S.K., Dwivedi Y.K., Metri B., Rana N.P. (eds) Re-imagining Diffusion and Adoption of Information Technology and Systems: A Continuing Conversation. TDIT 2020. IFIP Advances in Information and Communication Technology, vol 618. Springer, Cham. https://doi.org/10.1007/978-3-030-64861-9_43

[20]  N. Sebe, M. S. Lew, I. Cohen, A. Garg and T. S. Huang, "Emotion recognition using a Cauchy Naive Bayes classifier," *Object recognition supported by user interaction for service robots*, Quebec City, Quebec, Canada, 2002, pp. 17-20 vol.1, doi: 10.1109/ICPR.2002.1044578.

[21]  Senthilnathan, Samithambe, Usefulness of Correlation Analysis (July 9, 2019). Available at SSRN: https://ssrn.com/abstract=3416918 or http://dx.doi.org/10.2139/ssrn.3416918.

[22]  R. Taylor, EdD, RDCS, "Interpretation of the Correlation Coefficient: A Basic Review", 1990 Research Article.

# INSTRUCTIONS TO CONTRIBUTORS

The main aim of International Journal of Artificial Intelligence and Expert Systems (IJAE) is to provide a platform to AI & Expert Systems (ES) scientists and professionals to share their research and report new advances in the field of AI and ES. IJAE is a refereed journal producing well-written original research articles and studies, high quality papers as well as state-of-the-art surveys related to AI and ES. By establishing an effective channel of communication between theoretical researchers and practitioners, IJAE provides necessary support to practitioners in the design and development of intelligent and expert systems, and the difficulties faced by the practitioners in using the theoretical results provide feedback to the theoreticians to revalidate their models. IJAE thus meets the demand of both theoretical and applied researchers in artificial intelligence, soft computing and expert systems.

IJAE is a broad journal covering all branches of Artificial Intelligence and Expert Systems and its application in the topics including but not limited to technology & computing, fuzzy logic, expert systems, neural networks, reasoning and evolution, automatic control, mechatronics, robotics, web intelligence applications, heuristic and AI planning strategies and tools, computational theories of learning, intelligent system architectures.

To build its International reputation, we are disseminating the publication information through Google Books, Google Scholar, Open J Gate, ScientificCommons, Docstoc and many more. Our International Editors are working on establishing ISI listing and a good impact factor for IJAE.

The initial efforts helped to shape the editorial policy and to sharpen the focus of the journal. Started with Volume 10, 2021, IJAE appear with more focused issues related to artificial intelligence and expert systems studies. Besides normal publications, IJAE intend to organized special issues on more focused topics. Each special issue will have a designated editor (editors) – either member of the editorial board or another recognized specialist in the respective field.

We are open to contributions, proposals for any topic as well as for editors and reviewers. We understand that it is through the effort of volunteers that CSC Journals continues to grow and flourish.

## LIST OF TOPICS
The realm of International Journal of Artificial Intelligence and Expert Systems (IJAE) extends, but not limited, to the following:

- AI for Web Intelligence Applications
- AI Parallel Processing Tools

- AI Tools for Computer Vision and Speech Understand
- Application in VLSI Algorithms and Mobile Communication
- Case-based reasoning
- Derivative-free Optimization Algorithms
- Evolutionary and Swarm Algorithms
- Expert Systems Components
- Fuzzy Sets and logic

- Hybridization of Intelligent Models/algorithms
- Inference

- AI in Bioinformatics
- AI Tools for CAD and VLSI Analysis/Design/Testing
- AI Tools for Multimedia

- Automated Reasoning

- Data and Web Mining
- Emotional Intelligence
- Expert System Development Stages
- Expert-System Development Lifecycle
- Heuristic and AI Planning Strategies and Tools
- Image Understanding
- Integrated/Hybrid AI Approaches

- Intelligent Planning
- Intelligent System Architectures
- Knowledge-Based Systems
- Logic Programming
- Multi-agent Systems
- Neural Networks for AI
- Parallel and Distributed Realization of Intelligence
- Reasoning and Evolution of Knowledge Bases
- Rule-Based Systems
- Uncertainty

- Intelligent Search
- Knowledge Acquisition
- Knowledge-Based/Expert Systems
- Machine learning
- Neural Computing
- Object-Oriented Programming for AI
- Problem solving Methods

- Rough Sets
- Self-Healing and Autonomous Systems
- Visual/linguistic Perception

## CALL FOR PAPERS

# CONTACT INFORMATION

**Computer Science Journals Sdn BhD**

B-5-8 Plaza Mont Kiara, Mont Kiara

50480, Kuala Lumpur, MALAYSIA

Phone: 006 03 6204 5627

Fax:     006 03 6204 5628

Email: cscpress@cscjournals.org