# Prevention Pre-Violence in E-Labs with Machine Learning : PVE

Ahmet Furkan Aydogan
axa184@shsu.edu
Dept. of Computer Science
Sam Houston State University
Huntsville 77341 TX - USA

Narasimha Shashidhar
karpoor@shsu.edu
Dept. of Computer Science
Sam Houston State University
Huntsville 77341 TX - USA

*Abstract*—**Digital Forensics continues to be one of the most needed areas of today. In particular, the difficulties experienced in the field of education during the pandemic period have carried the fields where training will be given to the digital parts. However, it also brought pre-pandemic concepts such as bully, violence and insult to cyber environments. Although the studies to improve the existing distance education applications generally focus on areas such as quality image and sound transfer, the ones that may cause bigger problems in the future are the acts called crime. This study aims to create a possible schema to easily complete an investigation that may arise in the future and obtain evidence by applying the Digital Forensics field, which was created to investigate cyber-crime in detail and deliver it to judicial authorities, to distance education applications called e-lab. In addition, it can be applied to living systems to prevent the aforementioned criminal elements. While the study is performing itself, it focuses on machine learning and natural language processing, and it is seen that it has achieved more than 90% success in small-scale experiments.**

**Keywords:** Digital Forensics , E-Labs , Machine Learning , Natural Language Processing.

## I. INTRODUCTION

Although Electronic Discovery (e-Discovery) is not a popular phenomenon in the IT world, it continues to be an important sector and an emerging area for research. The concept of metadata, which appeared on computers in the early 1960s, would give rise to many innovations. In particular, there would be different challenges in protecting and listing metadata, followed by both legal and IT-based steps in filing and listing. In 1970, the concept of data emerged in accordance with rule 34 of the Federal Rules of Civil Procedure (FRCP). In addition, the said rule decided that the data to be discovered should only be made on devices posing a threat to the legal authorities [1].

In the beginning of 1995, cybercrime would increase and interest in the e-Discovery sector would begin to grow equally. At this time, the European Union (EU) decided to realize the importance of the digital data of its citizens and make innovations. The next year, the EU, which presents a Data Protection Directive, publishes confidentiality and progress reports on the bill under the heading "Processing".The important part was that the title "Processing" was based on a legal basis [2].

In the 2000s, the Department of Commerce in the United States accelerated its work and submitted the bill, Safe Harbor, to the government. Safe Harbor handled corporate confidentiality and contained a number of measures. The similarity with the Data Protection Directive developed by the EU was in fact the main reference. In 2015, the Safe Harbor bill was renamed as Privacy Shield [3].

In 2002, the so-called Enron Scandal revealed that the company's data was erased and the accounting records were cheated. The Enron issue would give rise to the Sarbanes-Oxley Act, to which the e-Discovery was constantly referred.The Sarbanes-Oxley Act contained articles by the US Congress to prevent companies like Enron from making false financial reports. The law, also known as the Corporate Responsibility Act, had strict criminal transactions and securities regulations [4].

Between 2003 and 2005, the e-Discovery standards came across an event that would be reorganized. In the so-called Zubulake case judge Shira Scheindlin of Southern District of New York would take steps to innovate in the field of e-Discovery. In particular, in the context of the case, new discourses like "companies should retain their data instead of destroying it because of the possibility of evidence in a proceeding" were mentioned. In addition, many new terms would be included in the e-Discovery such as electronic preservation, discovery, and legal hold notices. Finally, the Zubulake case would bring innovations to eliminate economic disparities between the company and plaintiffs. A further innovation in the field of pricing in the field of e-Discovery was to take place in 2007 between Mancia v. Mayflower Textile Services Co.. According to U.S. Magistrate Judge Paul Grimm of the District of Maryland's request , the lawyers of both parties must came together to discuss which data would be requested in which format. In addition, Grimm also requested the statistical data to determine how much cost the requests would cause [5].

Another innovation in the field of e-Discovery came between 2005 and 2006. First, in 2005, The Electronic Discovery Reference Model (EDRM) was created by two lawyers, Tom Gelbmann and George Scoha. EDRM, which is the first written model about of how Electronically Stored Information (ESI) steps should be, tried to set certain standards in the field of e-

Discovery. In 2006, the Rules 26 (a) (1), 26 (b) (2), 33 and 34, established by the FRCP, refer to the necessity of ESI steps. Finally, Rule 37 (f) defined ESI as a safe harbor for computer systems [6][7].
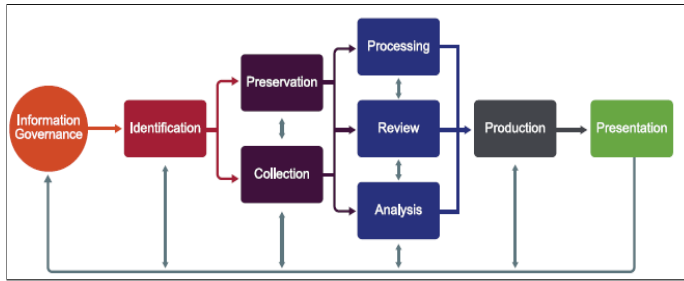


Fig. 1: Illustration of Steps of EDRM.

Another important title in the field of e-Discovery is Qualcomm Inc v. Broadcom Corp. Qualcomm's mistakes in the field of e-Discovery in the series of ongoing trials in 2007 brought significant penalties. First, the case, which involved troubled processes such as wrong numbering of 46,000 e-mails followed by document concealment, took a lean to against Qualcomm. As a result of the failure of external companies in the collection, numbering and filing of e-Discovery, $8.5 million was fined and patent rights were lost [8].

In 2008, the independent and non-profit legal research commission convened at the Sedona Conference. A team of lawyers commented on cost constraints in the field of e-Discovery. Specifically, in the ESI step, the analysis of the US government under great burden was put forward. In addition, it was mentioned that financial liabilities may cause disagreements between the parties and the deterrent power of cost.At the same year, the Federal Rule of Evidence (FRE) would have led to another innovation in the field of e-Discovery. According to FRE Rule 502, the defendant secures confidentiality and content that may contain sensitive data. In addition, FRE has restricted the use of Attorney-Client Privilege or Work-Product Doctrine to prevent disclosure. The said innovation is in the name of avoiding accidental waiver of the subject reported [9].

The year 2010 gained importance with the court decisions in the field of e-Discovery. First, the case between Victor Stanley Inc. and Creative Pipe Inc. was directed by Judge Grimm. After a while , case turned out to different way because the defendant blamed for failing to implement a legal hold and destroy the ESI after the lawsuit was filed. Meanwhile, the plaintiff, Victor Stanley Inc. , rights were automatically waived because refusal to cooperate. The court ordered the defendant to pay $ 1 million in attorneys' fees and two years in prison. In the another case between Pension Committee of the University of Montreal Pension Plan and Banc of America Securities , Judge Scheindlin took steps to stand against negligence. The judge stated that perfection is not important in e-Discovery but sustainability and balance of responsibility. Judge Scheindlin invited the plaintiff to take responsibility for the steps of the evidence to be collected.

In the end of 2010, there were incidents that led to the start of an event that would completely change the e-Discovery area. A study called Technology-assisted Review (TAR) would strengthen the relationship of e-Discovery with information technology. The article Document Categorization in Legal Electronic Discovery: Computer Classification vs. Manual Review, published by Anne Kershaw, Herbert L. Roitblat, and Patrick Oot, contained statistical data between the accuracy of the TAR system and the accuracy of manual research. In the beginning of 2011 , an article published in the Richmond Journal of Law and Technology by Maura R. Grossman and Gordon V. Cormack , would prove that the TAR system provides much more efficiency and accuracy than manual discoveries. After that moment , a great period was beginning for computer based e-Discovery [10][11][12].

In February 2012, in the case between Da Silva Moore and Publicis Groupe , U.S. Magistrate Judge Andrew J. Peck of the Southern District of New York approved the use of the TAR system. Judge Andrew, believing in the reliability of computer-generated data, would allow the TAR system to get approval from a judicial authority.Immediately afterwards, The American Bar Association (ABA) stated in their Rule 1.1 of a declaration no. 8 that technological approaches should bring ethical responsibilities.In 2013, in a case between Race Tires America and Hoosier Racing Corp. , FRCP 54 (d) (1) stated that rules 28 U.S.C. § 1920(4) , which is state that any material that is copied and charged necessarily obtained for use in the court , would not work in the field of e-Discovery.In addition, the court, which does not issue an appropriation to an expert to be designated for e-Discovery , has announced that it will only charge for documentation in TIFF format [13].

In 2015, the FRCP amended Rule 26 (b) (1) and Rule 37 (e) to address issues of proportionality and spoliation. In addition, FRCP Rule 26 (b) (1) stated that an amendment involving five factors would require the e-Discovery process to be carried out according to the needs of the case. Rule 37 (e) imposed severe criminal sanctions against irresponsibility and inadequate information. Another change in the beginning of 2015 was related to TAR application. The study Evaluation of Machine-Learning Protocols for Technology-Assisted Review in Electronic Discovery , published by Maura Grossman and Gordon Cormack, reported that the Continuous Active Learning (CAL) structure would yield superior results. For the implementation of the CAL structure, Machine Learning should be applied to the e-Discovery field and TAR 2.0 was released [14].

As mentioned above, e-Discovery has undergone radical changes. However, it is possible to examine the changes in three different categories. First of all, secrecy has been tried to be secured by both the court and institutions and it has gained an important place. In addition, a series of rules called General Data Protection Regulation (GDPR) was developed at a United Nations session in 2018 and the importance of data protection and confidentiality was once again addressed. Secondly, remuneration policies are in constant change. E-Discovery, requiring high-budget technological investigations,

led to changes in the wages of technological challenges as well as the fees of experts to conduct certification and research. Finally, with TAR 2.0, which pioneered technological research, e-Discovery was driven in a completely different direction. In particular, it is thought that privacy and pricing problems can be solved by using Machine Learning algorithms. Our article offers an application powered by Machine Learning, which aims to accelerate pre-evaluations in the field of e-Discovery as well as avoiding leakage of confidentiality and avoiding high fees for mentioned three main topics in e-Discovery. It can be used to prepare notifications and documents for both the main system and ongoing investigations by examining potential criminal elements that may arise in E-labs while performing transactions [15].

## II. WORKING SCHEMA OF PVE

The field of e-Discovery undoubtedly requires deep research. However, experts may want to ensure the consistency of the data to be investigated before conducting a deep investigation. Our method is to check the quality of the data to be investigated, or even to check whether the data is worthy of investigation. The Prevention Pre-Violence in E-Labs with Machine Learning (PVE), which can be called a trigger for investigation, uses the Natural Language Processing method as the main booster.In addition, PVE has an easy-to-understand interface for investigators.After the research, PVE provides different types of filing for recording. PVE will make a major contribution to the field of e-Discovery, which includes the recording of the investigator's information, conducting research between specific dates, taking important notes about the research, a changeable word pool, and a visual reportable result report.

### A. Preliminary Studies and Solution Schemes

This title contains the main information of PVE and can be described as a summary to explain exactly what the tool is doing.

*a) Understanding the Problems in the Field of e-Discovery:* As can be seen in the introduction part of the study, e-Discovery hosts different problems. In order to expand the research, a detailed examination of the problems experienced in the field of e-Discovery will be considered as the correct method.Data Variety, Increasing Data Volumes, Limited Human Resources, Declining Budges and Collaboration titles were identified as the main factors of the problems experienced in e-Discovery. The mentioned headings cause a complete loss of time and money. So much so that some e-Discovery charges can cost millions of dollars. In order to solve the problems, PVE provides a quick search on the text-based data such as e-mail, corporate messaging programs and SMS, which can be deemed appropriate by the court, in order to find the data that can constitute a criminal offense.
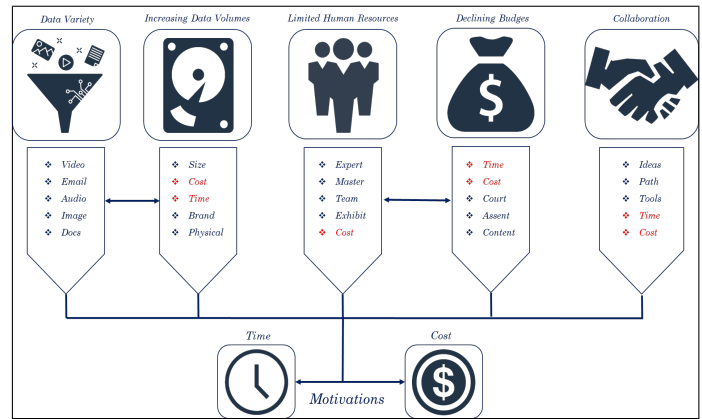


Fig. 2: Illustration of Problems of e-Discovery.

*b) Interface of PVE:*

PVE is a tool that has been taken care of in order to provide ease and intelligibility during use. The interface is enhanced with C # language, which emphasizes simplicity but does not limit its features in order to gain it.The interface starts by getting information about the person to be investigated. The requested information is as follows;

1) *Case Number.*
2) *Investigator Name.*
3) *Institution Name.*
4) *Department Name.*
5) *Investigation Date.*
6) *Special Notes About Investigation.*



Fig. 3: Illustration of PVE's Investiagator Record.

Following the registration of the person to be investigated, the PVE tool is transferred to the data section. In the data insertion section, the information, files and date ranges of the person to be investigated are requested. Then, four different legal grounds are used. Legal grounds have been prepared in the correct proportion with the contents to be determined within the investigation. For example, the title of the FRCrP contains a pool of words with correspondence such as terrorism, espionage, international crime, which can be deemed a

crime in government laws and the data added according to these words will be examined.In addition, the data section can cross-examine between legal bases. For example, if the defendant is considered to act against both the government and the detriment of one person, both pools of words can be used at the same time. In addition, after examining the rules for the corporation belonging to the companies, queries can be made as specified by contributing to the word pool. Finally, by creating a user-friendly interface for reporting the review, fortmats such as CSV, OCR, DOCX and PDF, which are most preferred by legal institutions, can be printed. In addition to the name, date and special notes of the investigator on the report to be printed, a detailed print of the examination may be taken. If requested, it is possible to report crime parts instead of the whole researched document. With the option specified, there will be a significant reduction in costs because it is more convenient to remove only the usable parts out of thousands of emails than to print them all.
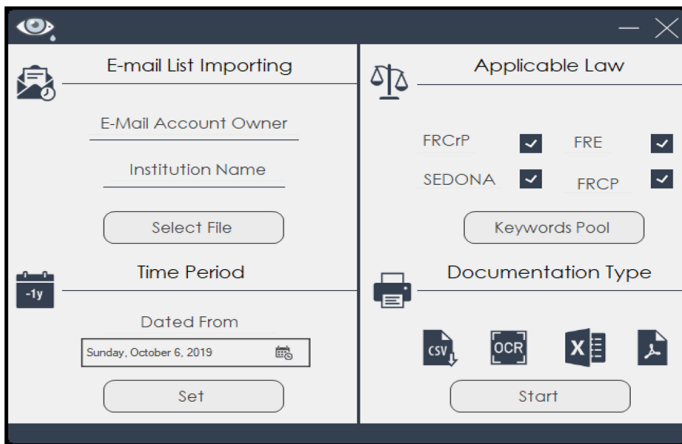


Fig. 4: Illustration of PVE's Data Analysis Settings Side.

Finally, PVE can provide visual reporting after the investigation. The picture, name and department information of defendant as well as word temperature map and pie chart are presented visually.
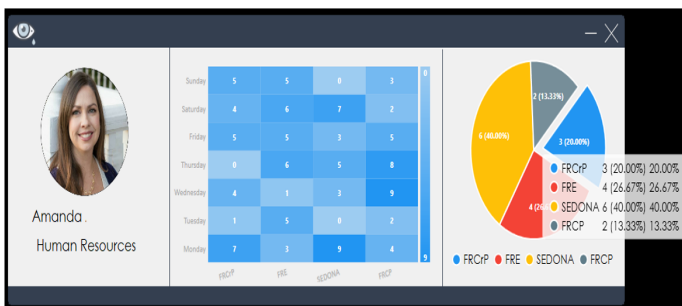


Fig. 5: Illustration of PVE's Visual Reporting Side.

c) *Possible Solution Schemes of PVE:*

The PVE tool is able to perform Machine Learning based operations. PVE uses C # for the user interface, and the Python programming language to perform background operations.

After entering the data to be reviewed, the data is sorted by date. The generated data set is transferred to the Python programming language. Utilities that were previously stored in a live system and built with Python libraries receive incoming data to perform their operations. The received data is ready for analysis after passing certain steps. The specific steps mentioned will ensure that the incoming data is ready for the Natural Language Process (NLP). Especially, Lemmatization, Stemming, Filtering Punctuation, Word Mapping, Keywords Mapping, Converting Words, Vocabulary Checking are used in PVE since they are suitable for NLP. The operating principle of the PVE tool is shown below with rough lines;
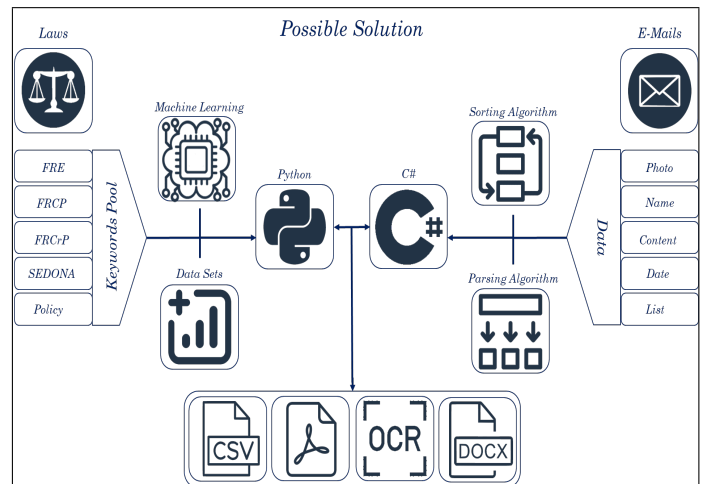


Fig. 6: The Operating Principle of the PVE.

The sections described describe what is PVE with its rough lines. The following headings will usually include details of background operations developed with the Python programming language.

B. *PVE's Detailed Work Mechanism.*

This section will provide an overview of key code snippets and strategies in the background of the PVE application. It should be noted that in NLP applications, it is important to first make the data understandable by the machine. Attention should be paid to features such as punctuation or capitalization that may lead to poor performance. The Python programming language is rich in NLP libraries and includes easy adaptations. The following subtitles have more detailed information.

a) *Possible Data-sets and Types:*

1) *Federal University of Sao Carlos (UFSCar) Corpus*

The SMS Spam Collection is a set of SMS tagged messages that have been collected for SMS Spam research. It contains one set of SMS messages in English of 5,574 messages, tagged according being ham (legitimate) or spam [16].

2) *rDany Chat*

157 chats 6300+ messages with a (fake) virtual companion.This bot have a personality ;

1) *Candid.*
2) *True.*
3) *Fun.*
4) *Optimistic.*
5) *Empathetic.*
6) *Gender Neutral.*
7) *Likes Art.*

3) *A semi-comprehensive List of Profanity in English.*

Messaging document for a profanity in English. It contains 2850 words.

4) *Google's Profanity Words*

Full List of Bad Words and Top Swear Words Banned by Google. It contains 451 words.

5) *Carnegie Mellon School of Computer Science's Offensive Words*

A list of 1,300+ English terms that could be found offensive. The list contains some words that many people won't find offensive [17].

6) *Data and Code for the Study of Bullying*

Version 3.0: bullying V3.0.zip (size 534950, released in June 2015). 7321 tweets with tweet ID, bullying, author role, teasing, type, form, and emotion labels.This version was described in: Junming Sui. Understanding and Fighting Bullying with Machine Learning. PhD thesis, Department of Computer Sciences, University of Wisconsin-Madison. 2015.(Archived version) bullying V2.0.zip (size 217680, released in September 2014). 1762 tweets with tweet ID, bullying, author role, and teasing labels.(Archived version) bullying V1.zip (size 19141, released in April 2012). Same tweets as in V2.0 but without tweet IDs.Versions 1.0 and 2.0 of this data set were introduced in the paper: Jun-Ming Xu, Kwang-Sung Jun, Xiaojin Zhu, and Amy Bellmore. Learning from bullying traces in social media. Proceedings of NAACL HLT 2012 [18].

7) *WhatsApp Dataset on Cyberbullying*

WhatsApp dataset to study cyberbullying among Italian students aged 12-13 in the context of the CREEP EIT project.The corpus of Whatsapp chats is made of 14,600 tokens divided in 10 chats. All the chats have been annotated by two annotators using the CAT web-based tool following the same guidelines [18].

It should be noted that the aforementioned data sets can note used at the same time. Specifically, sets 1 and 2 could use for character analysis and filtering. In addition, the data sets 3,4,and 5 were used for a general word pool. Finally, data sets 6 and 7 could use for decomposition for FRCP, Policy,

and general guilt detection. The following sub-headings describe the detailed use of data sets.

*b) Normalization of Data-sets with Natural Language Toolkit:*

In this area, the use of the Natural Language ToolKit (NLTK) library, which is used to perform Natural Language Processing operations, with python, and the meaning of the features such as classification, tokenization, stemming, tagging, parsing, and semantic reasoning, which belong to NLTK, are explained.

1) *Tokenization with NLTK*

Tokenization process, as a feature in NLTK, allows easier processing of sentences in the separation into vocables and in the later stages. In particular, it provides great advantage in punctuation and some other cleaning operations and helps to improve performance.

2) *Stemming Words with NLTK*

Stemming algorithms play a helpful role in separating the words in the sentence from time adverbs and making them simpler. In particular, it is an important process in order to facilitate the operations of machine learning algorithms and to get more precise results.

3) *Lemmatization Words with NLTK*

Lemmatization processes are used to bring the words whose roots are determined back to their main state. Uses the most useful and meaningful synonym for the root word while performing its operations.
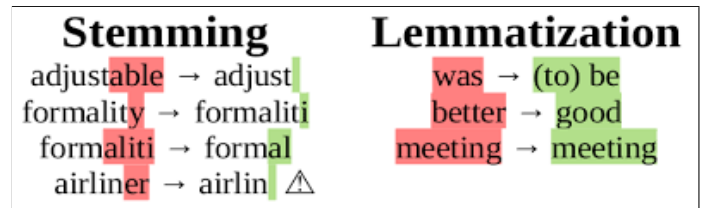


Fig. 7: Illustration of Differences between Stemming and Lemmatization.

4) *Stop-Words with NLTK*

Stop words operations are used to eliminate the parts in a sentence that can be considered unnecessary. It is a known fact that sentences free of grammatical features such as preposition and post-position in English provide a great performance increase in machine learning algorithms.

6) *Elimination of Punctuation, Capital Letters ,and Spell Checking*

Another important issue for NLP is the purification of words, of course, punctuation marks, capitalization and incorrect word usage. With a fast and effective algorithm created, after the elimination of capital letters and punctuation marks are accomplished easily, using the library named SpellCheck, the most likely convergent word of the misused word can be selected.

## III. Experimental Results

Due to the detailed reporting of a system that can show all the features of the PVE application and the necessity to produce a hardware product on a large scale, the potential success rate with a small-scale application is shown in the following areas. Two different strategies have been tried using the dataset prepared by Zafarullah Mahmood for Preprocessed data for Toxic Comments Classification Challenge [19]. Within the mentioned dataset, wikipedia comments were included and it contains approximately 153.000 different comments. Comments include racism, cyberbullying, and sexism. In the scenario of the application, it is the examination of a text written by a student or user in an E-Lab environment and checking the violation capacity. By creating a wordlist belonging to the headings of racism, cyberbullying, and sexism, the sentences in the dataset were checked and when the words in the wordlist were matched, the sentence was defined by adding it to the dataset. While applying Machine Learning algorithms, 80% was used for training and the remaining 20% for testing.

As a first strategy, it is aimed to control the usage frequency of the words in the comments by creating keywords detection system. For the aforementioned strategy, using the library named Rake, the first three keywords of the cleared sentences were detected and approximately 450,000 different words were examined using Logistic Regression, SVM, and KNN algorithms. Although SVM was the most efficient one among the results, 95% success was achieved in the detection of non-violent words, while a low success rate of 40% was achieved in violent words.

In the second strategy, sentences were used as a whole, and 80% of the same dataset was used as training and 20% as a test. Machine Learning algorithms include Logistic Regression, SVM, and KNN algorithms in the previous strategy. Within the framework of the results obtained, the SVM classification algorithm has again had the highest success rate. While 97% success was achieved in detecting non-violent sentences, 85% success was achieved in detecting violent sentences.

```
SVM
                precision    recall   f1-score   support

          0       0.95       0.99      0.97        1816
          1       0.95       0.77      0.85         382

   accuracy                            0.95        2198
  macro avg       0.95       0.88      0.91        2198
weighted avg       0.95       0.95      0.95        2198
```

Fig. 8: Illustration of SVM Classifier Results.

## IV. Conclusion

The need for digital forensics has increased gradually during the pandemic period. Cyber violence rates are increasing in parallel, especially with distance education, home offices and many work areas becoming online. The Prevention Pre-Violence in E-Labs with MachineLearning (PVE) application enables a machine learning-based digital forencisc research to be conducted in the aforementioned cyber workspaces and E-Lab environments. Although digital forencisc investigations with e-discovery cause large budgets, the budget can be reduced with the PVE application. This study includes an experiment that can achieve 97% success with a small sample, while explaining the schemes and the course of how to apply a machine learning-based digital forencisc tool. In the later stages, it can be used to measure the potential violation rate in the profiles of users in the same system, sarcastic interpretations, live systems and E-lab environments.

## References

[1] D. W. Oard, J. R. Baron, B. Hedin, D. D. Lewis, and S. Tomlinson, "Evaluation of information retrieval for e-discovery," *Artificial Intelligence and Law*, vol. 18, no. 4, pp. 347–386, 2010.

[2] G. Pailli, "Can europe learn from us e-discovery?," in *Intellectual Property Forum: Journal of the Intellectual Property Society of Australia and New Zealand*, no. 96, pp. 44–54, 2014.

[3] T. Y. Allman, "The case for a preservation safe harbor in requests for e-discovery," *Def. Counsel. J.*, vol. 70, p. 417, 2003.

[4] Y. Li, "The case analysis of the scandal of enron," *International Journal of business and management*, vol. 5, no. 10, p. 37, 2010.

[5] R. C. Losey, "Mancia v. mayflower begins a pilgrimage to the new world of cooperation," in *Sedona Conf. J.*, vol. 10, p. 377, HeinOnline, 2009.

[6] J. G. Conrad, "E-discovery revisited: the need for artificial intelligence beyond information retrieval," *Artificial Intelligence and Law*, vol. 18, no. 4, pp. 321–345, 2010.

[7] E. Casey, *Handbook of digital forensics and investigation*. Academic Press, 2009.

[8] R. Schermerhorn, "Broadcam corp. v. qualcomm inc. 543 f. 3d 683 (fed. cir. 2008)," *DePaul J. Art Tech. & Intell. Prop. L.*, vol. 19, p. 203, 2008.

[9] J. M. Barkett, "Evidence rule 502: The solution to the privilege-protection puzzle in the digital era," *Fordham L. Rev.*, vol. 81, p. 1589, 2012.

[10] G. V. Cormack and M. R. Grossman, "Engineering quality and reliability in technology-assisted review," in *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pp. 75–84, 2016.

[11] H. L. Roitblat, A. Kershaw, and P. Oot, "Document categorization in legal electronic discovery: computer classification vs. manual review," *Journal of the American Society for Information Science and Technology*, vol. 61, no. 1, pp. 70–80, 2010.

[12] L. Carroll, "The grossman-cormack glossary of technology-assisted review," *Federal Courts Law Review*, vol. 7, no. 1, 2013.

[13] B. B. Borden and J. R. Baron, "Finding the signal in the noise: Information governance, analytics, and the future of legal practice," *Richmond Journal of Law & Technology*, vol. 20, no. 2, p. 7, 2014.

[14] G. V. Cormack and M. R. Grossman, "Scalability of continuous active learning for reliable high-recall text classification," in *Proceedings of the 25th ACM international on conference on information and knowledge management*, pp. 1039–1048, 2016.

[15] P. Voigt and A. Von dem Bussche, "The eu general data protection regulation (gdpr)," *A Practical Guide, 1st Ed., Cham: Springer International Publishing*, vol. 10, p. 3152676, 2017.

[16] J. M. G. Hidalgo, T. A. Almeida, and A. Yamakami, "On the validity of a new sms spam collection," in *2012 11th International Conference on Machine Learning and Applications*, vol. 2, pp. 240–245, IEEE, 2012.

[17] B. Goel and R. Sharma, "Usf at semeval-2019 task 6: Offensive language detection using lstm with word embeddings," in *Proceedings of the 13th International Workshop on Semantic Evaluation*, pp. 796–800, 2019.

[18] J. Sui, *Understanding and fighting bullying with machine learning*. PhD thesis, The University of Wisconsin-Madison, 2015.

[19] Z. Mahmood, "Preprocessed data for toxic comments classification challenge," Kaggle, 2018.