# DETECTION OF POSSIBLE ILLICIT MESSAGES USING NATURAL LANGUAGE PROCESSING AND COMPUTERVISION ON TWITTER AND LINKED WEBSITES

**Dr.D.Maria manuel vianny[1], T.Ramya Priya[2], Achina Manikanta[3], Banala Rajashekhar[4], Chinnabathni Nishanth[5]**

[1]Professor, Department of CSE, Sri Indu Institute of Engineering & Technology, Hyderabad

[2]Assistant Professor, Department of CSE, Sri Indu Institute of Engineering & Technology, Hyderabad

[3,4,5] IV[th] Btech Student, Department of CSE, Sri Indu Institute of Engineering & Technology, Hyderabad

## ABSTRACT

Human trafficking is a global problem that strips away the dignity of millions of victims. Currently, social networks are used to spread this crime through the online environment by using covert messages that serve to promote these illegal services. In this context, since law enforcement resources are limited, it is vital to automatically detect messages that may be related to this crime and could also serve as clues. In this paper, we identify Twitter messages that could promote these illegal services and exploit minors by using natural language processing. The images and the URLs found in suspicious messages were processed and classified by gender and age group, so it is possible to detect photographs of people under 14 years of age. The method that we used is as follows. First, tweets with hashtags related to minors are mined in real-time. These tweets are preprocessed to eliminate noise and misspelled words, and then the tweets are classified as suspicious or not. Moreover, geometric features of the face and torso are selected using Haar models. By applying Support Vector Machine (SVM) and Convolutional Neural Network (CNN), we are able to recognize gender and age group, taking into account torso information and its proportional relationship with the head, or even when the face details are blurred. As a result, using the SVM model with only torso features has a higher performance than CNN.

Keywords: Support Vector Machine (SVM), Convolutional Neural

# 1.INTRODUCTION

Initially the websites were isolated and just placed for reading since the user could not truly interact with the web. However, from the innovation and arrival of web 2.0, there was a revolutionary and radical change since the user stopped being a simple spectator and became an active individual in social networks such as Facebook, Twitter, Instagram, among others.

Unfortunately, a door has also been opened for illegal businesses such as human trafficking, where some countries, such as Latin American countries, have the highest rates of smuggling of people, especially children and adolescents under 14 years old. It is important to note that the average age of consent is 14 years old in Latin American countries, so if underage people are used for illicit services are directly considered victims of human trafficking. Currently, in Twitter, it is possible to find websites that offer escort or similar services where young girls are promoted for the consumption of ''customers.'' These girls are generally abused physically, psychologically, and sexually.

In, the authors use computer vision algorithms to predict age with an approximate accuracy of 86.64%. In, SVM and CNN classification models are used to define the gender of a person. To the best of our knowledge, there are no works that consider characteristics of the upper body (upper torso) in the images to classify age groups The present work has two phases. In the first stage, natural language processing techniques are used in order to identify messages on Twitter that promote illicit services provided by minors. In the second phase, from the websites categorized as suspects, images are extracted in order to perform image processing and gender recognition of two age groups: over 14 years and under or equal to 14 years old. For this recognition, not only the characteristics of the torso but also the facial features were used.

## 1.1 MOTIVATION

The motivation for detecting possible illicit messages using NLP and CNN is to develop an automated system that can detect and flag messages that contain illicit content, such as hate speech, threats, or illegal activity. The objective is to create a more efficient and effective method for identifying illicit messages than traditional manual methods, which are often time-consuming and subjective.

## 1.2 PROBLEM STATEMENT

The problem statement is that there is a growing need for accurate and efficient detection of illicit messages in various domains, such as social media, email, and messaging platforms. This is especially important for businesses and organizations that want to ensure the safety of their employees and customers and comply with legal and ethical standards. However, identifying and categorizing messages manually can be challenging due to the sheer volume of messages and the subjectivity of human interpretation.

## 1.3 OBJECTIVE OF THE PROJECT

The objective of the project is to develop an automated system that can accurately and

efficiently identify and categorize illicit messages. This involves training a CNN model on a large dataset of labeled messages, developing a pre-processing pipeline for the messages, and evaluating the performance of the model on a test dataset. The system can then be used in real-time to detect and flag messages containing illicit content. The ultimate goal is to improve the safety and security of online communication by providing a more reliable and efficient method for detecting illicit messages.

# 2. LITERATURE SURVEY

## 2.1 A NON-PARAMETRIC LEARNING APPROACH TO IDENTIFY ONLINE HUMAN TRAFFICKING

Human trafficking is among the most challenging law enforcement problems which demands persistent fight against from all over the globe. In this study, we leverage readily available data from the website "Backpage"- used for classified advertisement- to discern potential patterns of human trafficking activities which manifest online and identify most.

## 2.2 A NEW ALGORITHM FOR AGE RECOGNITION FROM FACIAL IMAGES

This work provides an age-group recognition algorithm based on frontal face pictures. The algorithm's four primary stages include pre-processing, extraction of facial features using a unique geometric feature- based technique, analysis of facial features, and age classification. To utilise the approach, we require a face image database that provides age-related information about the people in the photographs. The unavailability of such a database drove us to create our own, which we termed the Iranian face database (IFDB). The IFDB holds digital photos of people ranging in age from one year to eighty-five years. Following pre-processing, the essential features of a database may be consistently determined. Finally, a neural network is used to classify face parameters such as wrinkle density and ratios. According to the results of an experiment, the system accurately classifies the age range with a precision of 86.64 percent.

## 2.4 EXISTING SYSTEM

As there is no staff available in unmanned restaurants, it is difficult for the restaurant management to estimate how the concept and the food is experienced by the customers. Existing rating systems, such as Google and TripAdvisor, only partially solve this problem, as they only cover a part of the customer's opinions. These rating systems are only used by a subset of the customers who rate the restaurant on independent rating platforms on their own initiative. This applies mainly to customers who experience their visit as very positive or negative.

## 2.4.1 DISADVATAGES OF EXISTING SYSTEM

Although there are previous tweet filtering and image classification works to detect illicit messages, most of them use

natural language processing methods or computer vision techniques.

## 2.5 PROPOSED SYSTEM

In order to solve the above problem, all customers must be motivated to give a rating. This paper introduces an approach for a restaurant rating system that asks every customer for a rating after their visit to increase the number of ratings as much as possible. This system can be used unmanned restaurants; the scoring system is based on facial expression detection using pretrained convolutional neural network (CNN) models. It allows the customer to rate the food by taking or capturing a picture of his face that reflects the corresponding feelings. Compared to text-based rating system, there is much less information and no individual experience reports collected. However, this simple fast and playful rating system should give a wider range of opinions about the experiences of the customers with the restaurant concept.

In this paper author first crawling twitter by using words like Lolita, escort and many more and then extracted tweets will go for cleaning to remove special symbols and stop words (words such as the, where, and, an, are etc.) and then tweets will be analyze to extracts words such as VERBS and ADJECTIVE and these words may contain important subjects or suspicious words used by HUMAN TRAFFICKERS (the suspicious words can be chicken soup, girls, penguin and many more. Clean tweets will be given input to

SVM and Naïve Bayes classifier to detect suspicious words.

If any tweet contains suspicious words then that tweet website will be scanned for images and each image will be processed through SVM HAARCASCADE classifier to detect face from that image and same algorithm will be used to detect upper body and both resultant images will be input to CNN (Convolution Neural Networks) classifier which will detect or predict AGE and GENDER from the resultant images. In this paper we are detecting gender as MALE and FEMALE and AGE will predicted with two classes as UNDER 14 Years or OVER 14 Years.

### 2.5.1 ADVANTAGES OF PROPOSED SYSTEM

Using predictive models, such as Vector Support Machine (SVM) and Convolutional Neural Networks (CNN), the image classification process is done through a training phase and a testing phase.

## 2.6 FEASIBILITY STUDY

The feasibility of the project is analyzed in this phase and business proposal is put forth with a very general plan for the project and some cost estimates. During system analysis the feasibility study of the proposed system is to be carried out. This is to ensure that the proposed system is not a burden to the company. For feasibility analysis, some understanding of the major requirements for the system is essential.

**Three key considerations involved in the feasibility analysis are,**

- ➢ **ECONOMIC FEASIBILITY**
- ➢ **TECHNICAL FEASIBILITY**
- ➢ **SOCIAL FEASIBILITY**

## 2.6.1 ECONOMIC FEASIBILITY

This study is carried out to check the economic impact that the system will have on the organization. The amount of fund that the company can pour into the research and development of the system is limited. The expenditures must be justified. Thus the developed system as well within the budget and this was achieved because most of the technologies used are freely available. Only the customized products had to be purchased.

## 2.6.2 TECHNICAL FEASIBILITY

This study is carried out to check the technical feasibility, that is, the technical requirements of the system. Any system developed must not have a high demand on the available technical resources. This will lead to high demands on the available technical resources. This will lead to high demands being placed on the client. The developed system must have a modest requirement, as only minimal or null changes are required for implementing this system.

## 2.6.3 SOCIAL FEASIBILITY

The aspect of study is to check the level of acceptance of the system by the user. This includes the process of training the user to use the system efficiently. The user must not feel threatened by the system, instead must accept it as a necessity. The level of acceptance by the users solely depends on the methods that are employed to educate the user about the system and to make him familiar with it. His level of confidence must be raised so that he is also able to make some constructive criticism, which is welcomed, as he is the final user of the system.

## 2.7 FEATURES OF THE PROJECT

- ➢ Computer vision offers the ability to sense surroundings and process the information it's taken in.
- ➢ NLP enables the understanding of spoken or written language and knowing which words to string together to communicate a prescribed message, much the same way as humans do.
- ➢ Integrating NLP and computer vision enables the design of assistive technology solutions for people who are deaf, converting their sign language into visuals or text.

## 2.8 TECHNOLOGIES USED FOR IMPLEMENTATION

### 2.8.1 PYTHON

Python was the major technology used for the implementation of machine learning concepts the reason being that there are numerous inbuilt methods in the form of packaged libraries present in python. Following are prominent libraries/tools we used in our project.

### 2.8.2 TENSORFLOW

TensorFlow is a free and open-source software library for machine learning and artificial-intelligence. It can be used across a range of tasks but has a particular focus on training and inference of deep neural

networks. TensorFlow serves as a core platform and library for machine learning. TensorFlow's APIs use Keras to allow users to make their own machine learning models.

### 2.8.3 SCIPY

SciPy is a free and open-source Python library used for scientific computing and technical computing. SciPy contains modules for optimization, linear algebra, integration, interpolation, special functions, FFT, signal and image processing, ODE solvers and other tasks common in science and engineering.

### 2.8.4 SCIKIT-LEARN

Scikit-learn provides a range of supervised and unsupervised learning algorithms via a consistent interface in Python. It is licensed under a permissive simplified BSD license and is distributed under many Linux distributions, encouraging academic and commercial use.

### 2.8.5 JUPYTER NOTEBOOK

The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations, and narrative text. It includes data cleaning and transformation, numerical simulation, statistical modeling, data visualization, machine learning. and much more.

### 2.8.6 PANDAS

Pandas is an open-source library that is made mainly for working with relational or labeled data both easily and intuitively. It provides various data structures and operations for manipulating numerical data and time series. This library is built on top of the NumPy library. Pandas is fast and it has high performance & productivity for users.

### 2.8.7 FLASK FRAMEWORK

Flask is a micro web framework. It is a Python module that lets you develop web applications easily. Flask is easy to learn due to smaller size and first in choice in case an application has to be developed quickly without going into much details of the framework. This Framework is used for HTML generation that a developer could possibly need to build an application. 'pip install Flask' command is used to install Flask framework in Python.

### 2.8.8 MATPLOTLIB

Matplotlib is an amazing visualization library in Python for 2D plots of arrays. Matplotlib is a multi-platform data visualization library built on NumPy arrays and designed to work with the broader SciPy stack. One of the greatest benefits of visualization is that it allows us visual access to huge amounts of data in easily digestible visuals. Matplotlib consists of several plots like line, bar, scatter, histogram etc**.**

### 2.8.9 HTML

An acronym for Hyper Text Markup Language it is a standard markup language that is used for designing and creating documents that would be displayed on any web browser. It can be further supported by technologies like

Cascading Style Sheets and JavaScript as a scripting language.

## 2.8.10 CSS

It stands for Cascading Style Sheets which is a style sheet language that defines the presentation of any document written using a markup language like HTML.

## 2.8.11 DJANGO

Django is a high-level Python Web framework that encourages rapid development and clean, pragmatic design. Built by experienced developers, it takes care of much of the hassle of Web development, so you can focus on writing your app without needing to reinvent the wheel. It's free and open source. Django's primary goal is to ease the creation of complex, database-driven websites.

# 3. ANALYSIS

## 3.1 INTRODUCTION

In recent years many criminal organizations advertise these ''sexual services'' using social networks hiding their illegal activity with seemingly innocuous terms such as ''chicken soup'' to refer to child pornography. Websites and social networks are used to extend this crime to the online environment, where covert advertising and messages are used to promote illegal services to exploit people who are victims of this crime, mainly minors. Although there are previous tweet filtering and image classification works to detect illicit messages, most of them use natural language processing methods or computer vision techniques separately.

## 3.2 REQUIREMENT SPECIFICATION

Requirement Specification plays an important role to create quality software solution. Requirements are refined and analyzed to assess the clarity. Requirements are represented in a manner that ultimately leads to successful software implementation. The project involved analyzing the design of few applications so as to make the application more users friendly. To do so, it was really important to keep the navigations from one screen to the other well-ordered and at the same time reducing the amount of typing the user needs to do. In order to make the application more accessible, the browser version had to be chosen so that it is compatible with most of the browsers.

## 3.2.1 SOFTWARE AND HARDWARE REQUIREMENTS

### 3.2.1.1 SOFTWARE REQUIREMENTS

- ➤ OPERATING SYSTEM : Windows 11
  Pip                  : 2.7
- ➤ LANGUAGES          : Python
- ➤ BROWSER            : Chrome
- ➤ FORNTEND           : HTML,CSS
- ➤ BACKEND            : Django

### 3.2.1.2 HARDWARE REQUIREMENTS

- ➤ PROCESSOR   : DUAL CORE 2.4GHZ(i5 or i7)
- ➤ RAM         :        20GB
- ➤ HARD DISK   :        40GB

# 4. DESIGN

## 4.1 INTODUCTION

Systems design is the process or art of defining the architecture, components, modules, interfaces, and data for a system to satisfy specified requirements. One could see it as the application of systems theory to product development. There is some overlap and synergy with the disciplines of systems analysis, systems architecture and systems engineering.

### 4.1.1 INPUT DESIGN

The input design is the link between the information system and the user. It comprises the developing specification and procedures for data preparation and those steps are necessary to put transaction data in to a usable form for processing can be achieved by inspecting the computer to read data from a written or printed document or it can occur by having people keying the data directly into the system. The design of input focuses on controlling the amount of input required, controlling the errors, avoiding delay, avoiding extra steps and keeping the process simple. The input is designed in such a way so that it provides security and ease of use with retaining the privacy. Input Design considered the following things:

- ➢ What data should be given as input?
- ➢ How the data should be arranged or coded?
- ➢ The dialog to guide the operating personnel in providing input.

- ➢ Methods for preparing input validations and steps to follow when error occur.

#### 4.1.1.1 OBJECTIVES

1. Input Design is the process of converting a user-oriented description of the input into a computer-based system. This design is important to avoid errors in the data input process and show the correct direction to the management for getting correct information from the computerized system.

2. It is achieved by creating user-friendly screens for the data entry to handle large volume of data. The goal of designing input is to make data entry easier and to be free from errors. The data entry screen is designed in such a way that all the data manipulates can be performed. It also provides record viewing facilities. 3.When the data is entered it will check for its validity. Data can be entered with the help of screens. Appropriate messages are provided as when needed so that the user will not be in maize of instant. Thus the objective of input design is to create an input layout that is easy to follow.

### 4.1.2 OUTPUT DESIGN

A quality output is one, which meets the requirements of the end user and presents the information clearly. In any system results of processing are communicated to the users and to other system through outputs. In output design it is determined how the information is to be displaced for immediate need and also the hard copy output. It is the most important and direct source information to the user. Efficient

and intelligent output design improves the system's relationship to help user decision-making.

1. Designing computer output should proceed in an organized, well thought out manner; the right output must be developed while ensuring that each output element is designed so that people will find the system can use easily and effectively. When analysis design computer output, they should Identify the specific output that is needed to meet the requirements.

2. Select methods for presenting information.

3. Create document, report, or other formats that contain information produced by the system.

## 4.2 PROPOSED SYSTEM ARCHITECTURE

### 4.2.1 TWEET CLASSIFICATION

➢ Harvesting Process
➢ Cleaning and Preprocessing
➢ Text Normalization
➢ Features Extraction and Classification

### 4.2.2 IMAGE CLASSIFICATION

➢ Scrapping
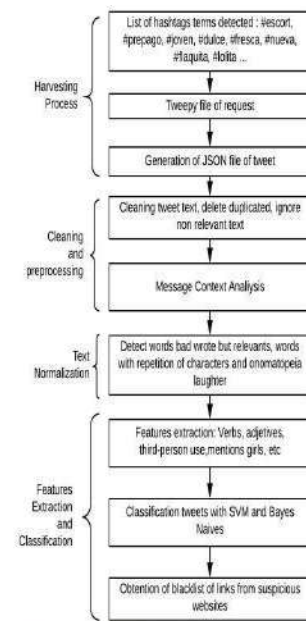➢ Feature Extraction
➢ Image Classification



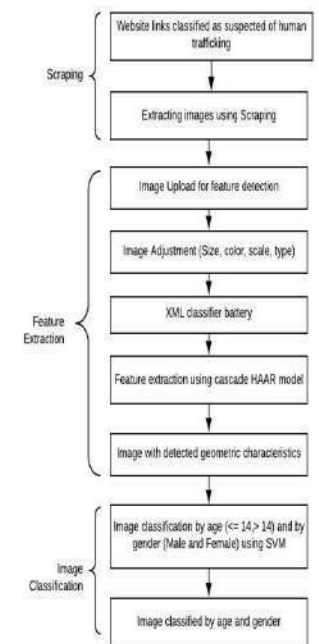FIGURE 1. Tweet classification based on natural language processing.



FIGURE 2. System overview of image processing.

## 4.3 UML DIAGRAMS

A UML diagram is a diagram based on the UML (Unified Modeling Language) with the purpose of visually representing a system along with its main actors, roles, actions, artifacts or classes, in order to better understand, alter, maintain, or document information about the system.

**Unified Modelling Language:** The Unified Modelling Language allows the software engineer to express an analysis model using the modelling notation that is governed by a set of syntactic semantic and pragmatic rules. A UML system is represented using five different views that describe the system from distinctly different perspective. Each view is defined by a set of diagrams, which is as follows

➢ User Model View
i.      This view represents the system from the user's perspective.

ii. The analysis representation describes a usage scenario from the end-user's perspective.

➢ Structural model view

i. In this model the data and functionality are arrived from inside the system.

ii. This model view models the static structures.

➢ Behavioral Model View

It represents the dynamic of behavioral as parts of the system, depicting the interactions of collection between various structural elements described in the user model and structural model view.

• Implementation Model View In this the structural and behavioral as parts of the system are represented as they are to be built.

• Environmental Model View In these the structural and behavioral aspects of the environment in which the system is to be implemented are represented.

## 4.3.1 USE CASE DIAGRAM

Use case diagrams are a set of use cases, actors, and their relationships. They represent the use case view of a system.

A use case represents a particular functionality of a system. Hence, use case diagram is used to describe the relationships among the functionalities and their internal/external controllers. These controllers are known as actors.
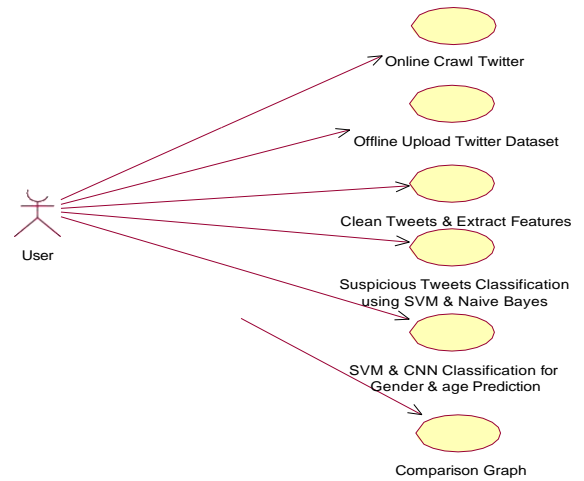


Figure 4.3.1- Use Case diagram

## 4.3.2 CLASS DIAGRAM

Class diagrams are the most common diagrams used in UML. Class diagram consists of classes, interfaces, associations, and collaboration. Class diagrams basically represent the object-oriented view of a system, which is static in nature.
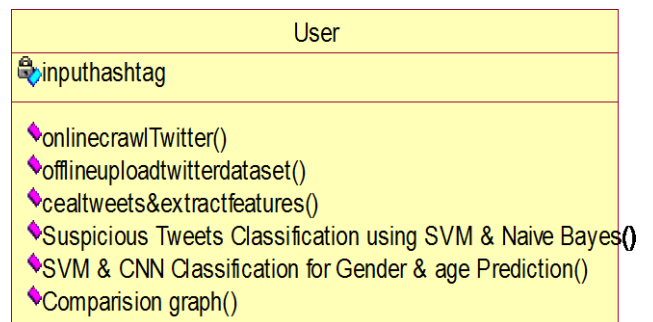


Figure 4.3.2 – Class diagram

## 4.3.3 SEQUENCE DIAGRAM

A sequence diagram is an interaction diagram. From the name, it is clear that the diagram deals with some sequences, which are the sequence of messages flowing from one object to another.
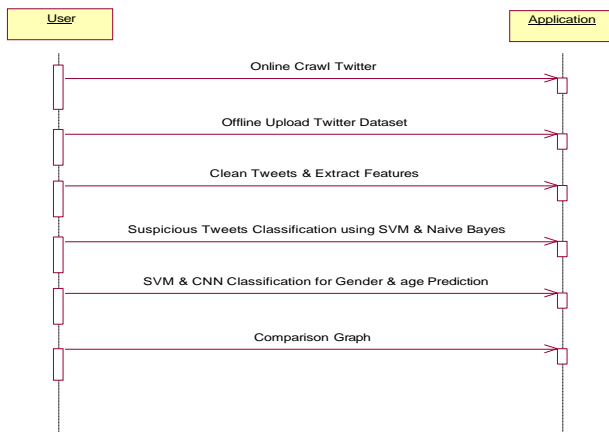
Figure 4.3.3 – Sequence diagram

## 4.3.4 COLLOBORATION DIAGRAM

The collaboration diagram is used to show the relationship between the objects in a system. Both the sequence and the collaboration diagrams represent the same information but differently. Instead of showing the flow of messages, it depicts the architecture of the object residing in the system as it is based on object-oriented programming. An object consists of several features. Multiple objects present in the system are connected to each other. The collaboration diagram, which is also known as a communication diagram, is used to portray the object's architecture in the system.
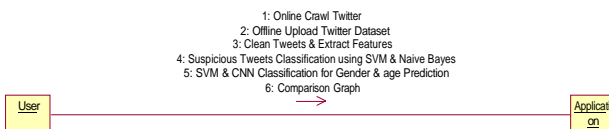


Figure 4.3.4 – Activity diagram

# 5. IMPLEMENTATION

## 5.1 IMPLEMENTATION

This section describes the process of tweets extraction, processing, and classification to determine if there are signs of human trafficking.

### 5.1.1 HARVESTING PROCESS

Initially, the data for each harvested day is stored on a JSON file that has information regarding the tweet post. The most relevant data includes the text of the tweet, user information, user mentions, associated URLs, and posted time. The capture process is shown in Figure. Spanish data is collected by executing a search request with the following hashtags: #escort, #prepago, #joven, #Dulce, #Fresca, #nueva, #lolita, and #flaquita. Hashtags were chosen as indicators of underage criteria. Tweets were used mined using the following criteria: mention of people from other countries if the tweets are written in the third person that shows that the Twitter user promotes the services of another person, or if the same user promotes the services of several people. To detect age, terms that indicate that people are underage victims were applied, such as the mention of a skinny young person or words from the jargon of pedophilia. A preliminary analysis of tweets and Facebook messages that were denounced as guilty of sex trafficking was conducted. The following words, in Spanish, are frequently used: joven (young), dulce (sweet), fresca (fresh), nueva (new), Lolita, and flaquita (skinny). Other words like Caldo de Pollo, club penguin, and cp are used for criminals as an abbreviation of child pornography, the hashtags #escort,

#prepago (prepaid), and the words mentioned before are chosen for this analysis. In Table 1, the number of tweet posts for each hashtag is summarized. For testing purposes, 100000 tweets were mined following the chosen words.

## 5.1.2 CLEANING AND PREPROCESSSING

All downloaded messages always contain at least one of the hashtags mentioned in Table 1, and these messages are stored locally in a JSON file. The information is processed and cleaned using a Python application, and tweets are deleted according to the following criteria:

| Hashtag | Number of occurrences |
|---|---|
| #escort | 45604 |
| #prepago | 15890 |
| #joven | 3456 |
| #dulce | 1256 |
| #fresca | 1456 |
| #nueva | 5743 |
| #flaquita | 6580 |
| #lolita | 867 |
| #penguin | 23980 |
| #caldodepollo | 45990 |
| #cp | 34562 |
| Twits with a URL link | 1765 |

Figure 5.1.2- Total tweets post by hashtags

➢ Tweets with certain characters and not standardized were removed in order to build a readable and more precise text.

➢ Repetitive tweets, a user may post the same tweet many times. Then, eliminating repetitions avoids redundant information; otherwise, we will create bias in the subsequent analysis.

➢ Tweet that does not contribute to the project.

For example, one of the words for the tweet filtering is ''young'' referring to children and adolescents, but if a user writes: ''Long live Quito! Young city!'', this message has another context. Therefore, it must be discarded. The cleaning of the data is essential to perform the classification, so digits, stop words, and special characters are removed. Moreover, duplicated advertisements or tweets out of context were eliminated, automatically. If a URL link was matched to a night club website or massage therapy site, the tweet was manually tagged.

## 5.1.3 TEXT NORMALIZATION

Twitter messages usually have much noise due to the shortness of the texts and because they are mostly generated using mobile devices. Besides, many tweets have incompleted, misspelled, or distorted words, so the performance of natural language processing is degraded. Consequently, in the preprocessing, it is necessary to apply methods of tweet standardization written in Spanish.

## 5.1.4 FEATURES EXTRACTION AND CLASSIFICATION

Features are defined based on some criteria related to the deception and cybercrime. These criteria consider young age as an indicator for the detection of victims. For the input characteristics, the reason why considered each one is explained in Table 3. With a Python program, syntax analysis is

done in order to evaluate the higher frequency of adjectives and verbs, which is valuable information since the deceptive message usually is very expressive.

The analysis of URL links to detect if they are night club or massage therapy site is made manually from a report produced by a program with a list of URLs, and the other features are obtained processing the corpus with some Python programs. Then, the data is loaded in the input feature file in order to feed the classifier. The characteristics of the twitter account and the nature of their messages were analyzed. Additionally, the tweets were tagged as ''suspicious'' when a message comes from accounts that were closed by Twitter due to complaints of child pornography content. A semi-supervised learning technique with Naïve Bayes and SVM algorithms was used in order to classify the tweets as ''suspicious'' or ''not suspicious'' of being related to sex trafficking. The performance of each classifier was evaluated based on average Precision (P), average Recall (R), and average F-Measure (F). Some algorithms were tested, and because SVM and Naïve Bayes presented a good performance and processing speed, they were chosen to classify the data using a semi-supervised approach. 10-fold cross-validation was used in order to evaluate the classifiers. This cross-validation divided the data randomly into ten sets. Each one was tested against the rest of the sets . The performance result was the average of all tests. As it is mentioned above, Precision, Recall, and F-Measure were used in order to evaluate classifiers' performance against evaluation using the ground truth established from the previous annotation.

## 5.3 ALGORITHM

## 5.3.1 SVM

Support Vector Machine (SVM) is a supervised machine learning algorithm used for both classification and regression. Though we say regression problems as well its best suited for classification. The objective of SVM algorithm is to find a hyperplane in an N-dimensional space that distinctly classifies the data points. The dimension of the hyperplane depends upon the number of features. If the number of input features is two, then the hyperplane is just a line. If the number of input features is three, then the hyperplane becomes a 2-D plane. It becomes difficult to imagine when the number of features exceeds three.

## 5.3.1.1 ADVANTAGES OF SVM

➢ Effective in high dimensional cases.

➢ Its memory efficient as it uses a subset of training points in the decision function called support vectors.

➢ Different kernel functions can be specified for the decision functions and its possible to specify custom kernels**.**

## 5.3.2 CNN

A Convolutional Neural Network (ConvNet/CNN) is a Deep Learning algorithm which can take in an input image, assign importance (learnable weights and biases) to various aspects/objects in the image and be

able to differentiate one from the other. CNNs are used for image classification and recognition because of its high accuracy. It was proposed by computer scientist Yann LeCun in the late 90s, when he was inspired from the human visual perception of recognizing things. A Convolutional neural network (CNN) is a neural network that has one or more convolutional layers and are used mainly for image processing, classification, segmentation and also for other auto correlated data.

# 6. TESTING, VALIDATION AND RESULTS

## 6.1 INRODUCTION

Testing is a process, which reveals errors in the program. It is the major quality measure employed during software development. During software development. During testing, the program is executed with a set of test cases and the output of the program for the test cases is evaluated to determine if the program is performing as it is expected to perform.

**Software testing** is the act of examining the artifacts and the behavior of the software under test by validation and verification. Software testing can also provide an objective, independent view of the software to allow the business to appreciate and understand the risks of software implementation.

A primary purpose of testing is to detect software failures so that defects may be discovered and corrected. Testing cannot establish that a product functions properly under all conditions, but only that it does not

function properly under specific conditions. The scope of software testing may include the examination of code as well as the execution of that code in various environments and conditions as well as examining the aspects of code: does it do what it is supposed to do and do what it needs to do. In the current culture of software development, a testing organization may be separate from the development team. There are various roles for testing team members. Information derived from software testing may be used to correct the process by which software is developed.

## 6.2 TESTING METHOLODIES

In order to make sure that the system does not have errors, the different levels of testing strategies that are applied at differing phases of software development are:

### 6.2.1 SYSTEM TESTING

System testing ensures that the entire integrated software system meets requirements. It tests a configuration to ensure known and predictable results. An example of system testing is the configuration oriented system integration test. System testing is based on process descriptions and flows, emphasizing pre-driven process links and integration points.

### 6.2.2 TYPES OF TESTING

### 6.2.2.1 UNIT TESTING

Unit testing involves the design of test cases that validate that the internal program logic is functioning properly, and that program inputs produce valid outputs. All decision

branches and internal code flow should be validated. It is the testing of individual software units of the application .it is done after the completion of an individual unit before integration. This is a structural testing, that relies on knowledge of its construction and is invasive. Unit tests perform basic tests at component level and test a specific business process, application, and/or system configuration.

## 6.2.2.2 INTEGRATION TESTING

Integration tests are designed to test integrated software components to determine if they actually run as one program. Testing is event driven and is more concerned with the basic outcome of screens or fields. Integration tests demonstrate that although the components were individually satisfaction, as shown by successfully unit testing, the combination of components is correct and consistent. Integration testing is specifically aimed at exposing the problems that arise from the combination of components.

## 6.2.2.3 FUNCTIONAL TESTING

Functional tests provide systematic demonstrations that functions tested are available as specified by the business and technical requirements, system documentation, and user manuals. Functional testing is centered on the following items:

Valid Input        : identified classes of valid input must be accepted.

Invalid Input        : identified classes of invalid input must be rejected.

Functions        : identified functions must be exercised.

Output        : identified classes of application outputs must be exercised.

Systems/Procedures  : interfacing systems or procedures must be invoked.

Organization and preparation of functional tests is focused on requirements, key functions, or special test cases. In addition, systematic coverage pertaining to identify Business process flows; data fields, predefined processes, and successive processes must be considered for testing. Before functional testing is complete, additional tests are identified and the effective value of current tests is determined.

**Each module can be tested using the following two Strategies:**

**Black Box Testing:**

Black Box Testing is testing the software without any knowledge of the inner workings, structure or language of the module being tested. Black box tests, as most other kinds of tests, must be written from a definitive source document, such as specification or requirements document, such as specification or requirements document. It is a testing in which the software under test is treated, as a black box .you cannot "see" into it. The test provides inputs and responds to outputs without considering how the software works.

**White Box testing:**

In this the test cases are generated on the logic of each module by drawing flow graphs of that module and logical decisions are tested on all the cases. It has been uses to generate the test cases in the following cases:

- ➢ Guarantee that all independent paths have been Executed.
- ➢ Execute all logical decisions on their true and false Sides.
- ➢ Execute all loops at their boundaries and within their operational bounds.
- ➢ Execute internal data structures to ensure their validity.

White Box Testing is a testing in which in which the software tester has knowledge of the inner workings, structure and language of the software, or at least its purpose.

## Test Approach:

Testing can be done in two ways:
- ➢ Bottom-up approach
- ➢ Top-down approach

**Bottom-up Approach:**

Testing can be performed starting from smallest and lowest level modules and proceeding one at a time. For each module in bottom up testing a short program executes the module and provides the needed data so that the module is asked to perform the way it will when embed within the larger system. When bottom level modules are tested attention turns to those on the next level that use the lower-level ones they are tested individually and then linked with the previously examined lower-level modules.

**Top-down approach:**

This type of testing starts from upper-level modules. Since the detailed activities usually performed in the lower-level routines are not provided stubs are written. A stub is a module shell called by upper-level module and that when reached properly will return a message to the calling module indicating that proper interaction occurred. No attempt is made to verify the correctness of the lower-level module.

## 6.3 VALIDATION

The system has been tested and implemented successfully and thus ensured that all the requirements as listed in the software requirements specifications are completely fulfilled. In case of erroneous input corresponding error messages are displayed.

Software validation checks that the software product satisfies or fits the intended use (high-level checking),

However, it is also possible to perform internal static tests to find out if the software meets the requirements specification but that falls into the scope of static verification because the software is not running.

## 6.3.1 TEST RESULTS

Test results are the outcome of the whole process of software testing life cycle. The results thus produced, offer an insight into the deliverables of a software project, significant in representing the status of the project to the stakeholders.

Reporting test execution results is very important part of testing, whenever test

execution cycle is complete, tester should make a complete test results report which includes the Test Pass/Fail status of the test cycle. If manual testing is done then the test pass/fail result should be captured in an excel sheet and if automation testing is done using automation tool then the HTML or XML reports should be provided to stakeholders as test deliverable. All the test cases mentioned above passed successfully.

# 7.CONCLUSION

## 7.1 CONCLUSION

Face recognition algorithms and machine learning models have been improved during the last years. For example, in the ILSVRC competition, an accuracy value of $90\% +-5\%$ was obtained. In these conditions, machine learning recognition can be similar to visual object recognition used by human beings. Many factors have a direct impact on image recognition, such as size, color, opacity, resolution, kind of image format, among others. Therefore, the results of image recognition and classification depend on the dataset quality.

In this work, we probed that satisfactory performance can be obtained using just geometric features of the torso and not only facial characteristics. For this paper, Haar filters combined with an SVM classifier were used for the extraction process of features, and then we classified the age group and gender with an SVM classifier. The obtained results were compared with the outcomes of a CNN algorithm.

To the best of our knowledge, this work is the first approach related to image classification without facial features but just the upper-body geometric characteristics. Currently, there is no similar research that takes into account only the upper body features of minors. Thus, the results of this paper can be applied to human trafficking, disappearance, kidnapping, among others. Moreover, the obtained information can be used by the police or other security institutions.

## 7.2 FUTURE ENHANCEMENT

Finally, future work includes:

1) the study of some characteristics related to ethnic and racial features,

2) to extend the proposal to extract geometric features of the entire body, another kind of images, or inclusive videos in different formats,

3) detection of medical issues by means the analysis of features extracted from torso images, legs, back, among other characteristics, and

4) the use of other algorithms or the applicability in other networks like Instagram.

# 8. BIBLIOGRAPHY

[1] B. Bangerter, S. Talwar, R. Arefi, and K. Stewart, ''Networks and devices for the 5G era,'' IEEE Commun. Mag., vol. 52, no. 2, pp. 90–96, Feb. 2014.

[2] F. Laczko, ''Data and research on human trafficking,'' Int. Migration, vol. 43, nos. 1–2, pp. 5– 16, Jan. 2005.

[3] M. Lee, ''Human trafficking and border control in the global south,'' in The Borders of Punishment: Migration, Citizenship, and Social Exclusion. Oxford, U.K.: Oxford Univ. Press, 2013, pp. 128– 149.

[4] E. Cockbain and E. R. Kleemans, ''Innovations in empirical research into human trafficking: Introduction to the special edition,'' Crime, Law Social Change, vol. 72, no. 1, pp. 1–7, Jul. 2019.

[5] R. Weitzer, ''Human trafficking and contemporary slavery,'' Annu. Rev. Sociol., vol. 41, pp. 223–242, Aug. 2015.

[6] T. S. Portal. (2018). Twitter: Number of Monthly Active Users 2010-2018. [Online]. Available: https://www.statista.com