# Classification with Explanation for Human Trafficking Networks

1st David Ing
*CRIL - CNRS, Université d'Artois*
Lens, France
ing@cril.fr

2nd Fabien Delorme
*CRIL - CNRS, Université d'Artois*
Lens, France
fabien.delorme@cnrs.fr

3rd Said Jabbour
*CRIL - CNRS, Université d'Artois*
Lens, France
jabbour@cril.fr

4th Nelly Robin
*Université de Paris - Campus Saint Germain des Prés*
Paris, France
nelly.robinsn@orange.fr

5th Lakhdar Sais
*CRIL - CNRS, Université d'Artois*
Lens, France
sais@cril.fr

*Abstract*—On a worldwide scale, an increasing number of victims of human trafficking were observed these last years, covering a majority of countries and territories. Among them, a large portion of women and girls are recruited primarily for sexual exploitation. United Nations Office on Drugs and Crime (UNODC) highlights the difficulties of access to justice which deprive victims of protection, a central issue behind our work. Our contribution is part of an emerging research trend, combining Artificial Intelligence (AI), Humanities and Social Sciences (HSS). It makes an original use of legal database to identify Human Trafficking Networks (HTNs), involving both sexual abuse victims and exploiters. First, a reformulation of the legal database as a numerical database is proposed, using new features expressing relationships between people involved in the same court case, likely to better reveal HTNs. Secondly, six machine learning algorithms, including Decision Tree, Random Forest, Gradient Boosting, Logistic Regression, Support Vector Machine (SVM) and K-Nearest Neighbors (KNN) are used to train on numerical database and learn to classify the input court case into one of the three classes: Not suspicious, Suspicious, or Probably suspicious. We in details discuss knowledge-based feature engineering, dataset balancing, parameters tuning, and best models selection. The comparative empirical evaluations between those classification algorithms have been conducted in order to highlights the relevance of our HTNs detection approach. To help the end-users, to better understand the displayed HTNs, for Decision Tree and Random Forest, we also provide explanations of why such court case can be classified. Those results were finally discussed with experts in the field of human trafficking, providing us with interesting feedback shedding light to this multidimensional form of modern-day slavery problem.

*Index Terms*—Machine Learning, Classification, Imbalanced data, Explainable AI, Humanities and Social Science, Human trafficking

## I. INTRODUCTION

Globally, 74,514 victims of trafficking were recorded between 2017 and 2018, distributed in 110 countries and territories [20][1]. In 2018, out of 10 identified victims worldwide, approximately five were adult women and two were girls, recruited primarily for sexual exploitation.

The International Organization for Migration (IOM) stresses that these statistics in themselves are deeply concerning, but the reality is far grimmer as it represents only a fraction of the true scale of human trafficking worldwide. Especially since "from the beginning of the Covid-19 pandemic, traffickers have been operating in greater clandestinity; therefore, it is more difficult to provide reliable estimates" [17]. The UNODC also highlights the difficulties of access to justice which deprive victims of protection. This situation reveals the social and political prominence of this issue and the importance of promoting the protection of victims. However, to effectively protect victims of trafficking, states must ensure that they are not held responsible for any offending activity as a direct result of their situation as a victim of trafficking, regardless of the seriousness of the offense committed , said a United Nations human rights expert on Tuesday (June 29, 2021). "Punishing a victim marks a break with the commitments made by states to recognize the priority of the rights of victims to assistance, protection and effective remedies", underlines Siobhàn Mullally, the special reporter on human trafficking in a report presented to the United Nations Human Rights Council (2021)[2]. The publication of this report coincides with a recent decision by the European Court of Human Rights which recalls that the punishment of a victim would be detrimental to his "physical, psychological and social recovery and could make him vulnerable to further trafficking in the future". In Senegal, the importance of the human trafficking issue is confirmed by the latest circular from the Ministry of Justice dated on November 5, 2021, addressed to all the courts of appeal, high and magistrates' courts. On the legislative and institutional level, with the ratification of the United Nations Convention against Organized Crime and its protocol on the fight against trafficking in persons, especially women and

---

[1]UNODC has been collecting data on trafficking in persons trends from official national judicial sources since 2003, and this edition of the Global Report on Trafficking in Persons draws on statistics spanning more than a decade, focusing is on the 2012-2014 period.

[2]ONU Info, 29 June 2021. https://news.un.org/fr/story/2021/06/1099222

children, the adoption of Law No 2005-06 of May 10, 2005[3] in this field and the creation of the National Unit for the Fight against Trafficking in Persons, Senegal has adopted a policy to fight against trafficking in persons.

This brief overview on the problem of trafficking in human beings highlights the main goal of this work, namely the design of a HTNs detection tool to protect human trafficking victims, especially women and children. To achieve this goal, we will rely on legal data, made available to us thanks to a host signed agreement between IRD[4] and ministry of justice of Senegal (see Section III-A); however, the goal is to produce a generic model in order to be able to extend it to all the courts of justice of the French-speaking countries of Economic Community of West African States (ECOWAS), which have the same legal standards and judicial practices as Senegal. To our knowledge, the legal database is used for the first time to address this important and complex issue. Our proposed framework is made of several processing steps as depicted in Figure 1(a). (1) The first step, from legal database, we aim to obtain a HTNs database containing potential HTNs (see step 1 in Section IV-A). (2) The HTNs database derived in (1) is annotated by an expert. (3) Then a reformulation of the annotated HTNs database is proposed (see step 2 in Section IV-A). (4) Figure 1(b), as we have an imbalanced dataset, resampling techniques are used with six machine learning classification algorithms, including Decision Tree, Random Forest, Gradient Boosting, Logistic Regression, Support Vector Machine (SVM) and K-Nearest Neighbors (KNN) to perform the comparative empirical evaluations. (5) Figure 1(c), for Decision Tree and Random Forest, the discovered networks are accompanied by explanations allowing end-users, the magistrates and human trafficking experts in our case, to understand and appreciate why and how they are classified. The results were finally discussed with experts in the field of human trafficking, providing us with interesting feedback.

The paper is organized as follows. We first discuss the related works about human trafficking (Section II). In section III, we provide some preliminaries, including the description of the legal database, the supervised classification algorithms that are used in this paper as well as the recent enhancements related to explainability of Decision Tree and Random Forest, then we describe the resampling approaches which were used to overcome the imbalanced classification problems. In Section IV, we describe how to reformulate the legal database to HTNs database, the experimental evaluations of our HTNs detection models, and the explanation of classification results. In section V, we discuss about the misclassifications made by those classifiers as well as the obtained results and their impacts, which have been analysed by experts before concluding.

## II. RELATED WORK

We observe recent interest on the problem of human trafficking detection mainly from social networks data. In [10], the

---

authors collect and process tweets to detect deception related to sex trafficking using predefined criteria as input features to machine learning classifiers with the task to classify the tweets as "suspicious" or "not-suspicious" of being related to sex trafficking. Considering escort websites a primary vehicle for selling the services of such trafficking victims, Wang et al. [21] proposed an ordinal regression neural network to identify escort ads that are likely linked to sex trafficking. Considering similar sources of data, escort websites, in [19] the authors achieve a major step in the automatic detection of advertisements suspected to pertain to human trafficking. They designed and trained a deep multi-modal model called the Human Trafficking Deep Network (HTDN). In [14] the authors present a novel unsupervised and scalable template matching algorithm for analyzing and detecting complex organizations operating on adult service websites. In [7], the authors present an approach using Natural Language Processing (NLP) to identify trafficking ads on websites offering sexual services through online advertisements. They propose a classifier by integrating multiple text feature sets, including the publicly available pretrained language model Bi-directional Encoder Representation from Transformers (BERT) [6]. As a summary, the few known approaches addressing the problem of human trafficking detection operate on escort websites or social networks, they are mainly based on machine learning, data mining or NLP techniques. From HSS side, most of the contributions mainly focused on the analysis of reports from public institutions, NGOs or lawyers.

We recall that human trafficking is a complex multidimensional problem for which we have incomplete data and limited knowledge about the actors and their interactions. Our paper aims to overcome these limitations, by mobilizing for the first time legal databases from Senegal, allowing to reveal HTNs and their structures, by mobilizing the machine learning classification techniques with explanations. The framework we propose brings a new dimension, the understanding of HTNs under the prism of the South.

## III. PRELIMINARIES

### A. Data Description

Let us briefly describe the legal database made available to us thanks to a host signed agreement between IRD and ministry of justice of Senegal. This database contains all criminal affairs or court cases (police and gendarmerie reports, complaints from individuals or public administrations) that come to the prosecution. They are recorded in the Register of Complaints. The follow-up given to each case by the public prosecutor is also recorded. Each criminal affair recorded in the Prosecutor's Registry (PR) includes the following attributes:

- Identification and location of the writer (*date, origin and place of the report or of the complaint*), Prosecutor's Registry Number (*PR Number*) or Judges Registry Number (*JR Number*), the person accused (*name, place of birth, nationality, address, profession, sex, age, family situation and number of children*).
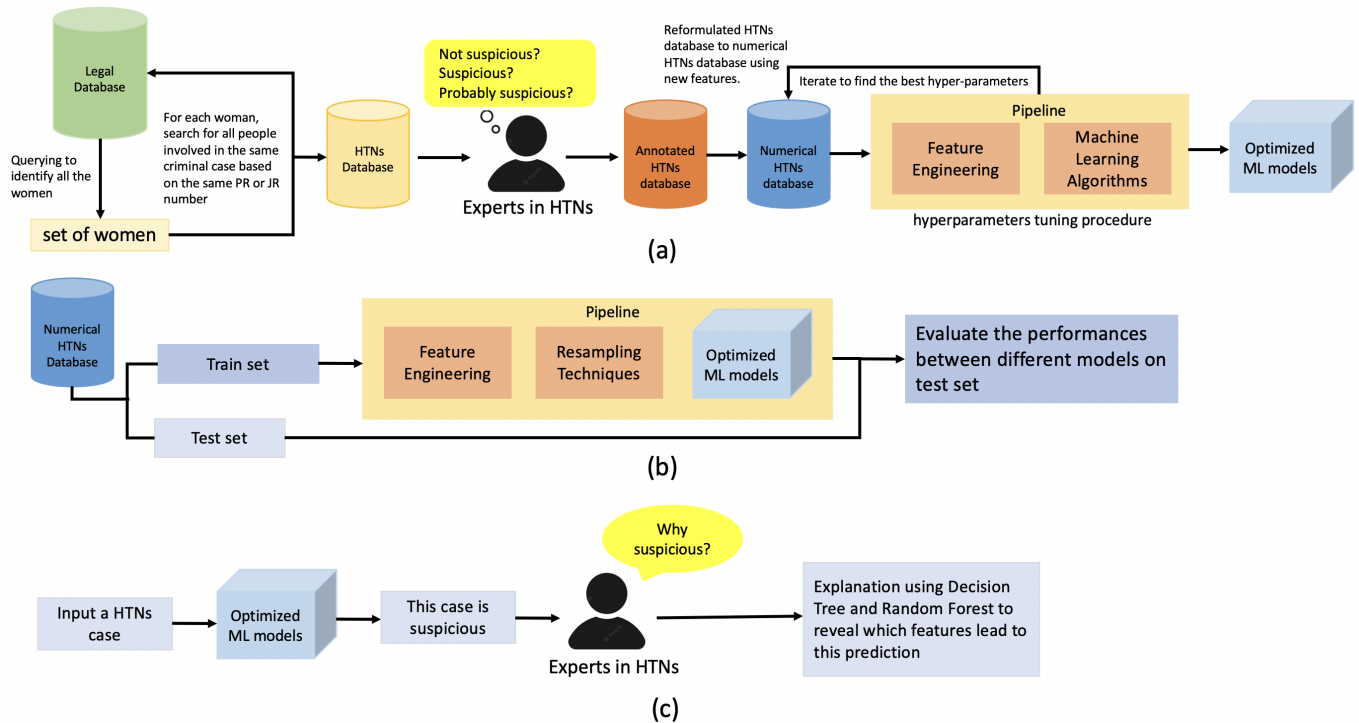
Fig. 1. The overall proposed framework for HTNs classification

- The qualification of the facts (*offenses committed*) which are reproached to the persons implicated (*lack of health and social booklet, pedophilia, pimping, prostitution, misappropriation of minors, soliciting, vagrancy, illegal residence, infanticide, swindling, possession and use of Indian hemp*, etc.), as well as the articles of laws which repress these facts.
- The follow-up reserved for the cases and the criminal situation of the persons implicated (*type of settlement: Flagrant offence, direct summons and introductory indictment*). The *court seized* (*TR, TE, SP*). The decision of the public prosecutor is also mentioned (*criminal situation: not in custody, warrant of deposition, temporary custody order for minors, status, duration of the sentence, etc.*).

It is important to note that several people can be involved in the same criminal affair. In such a case, they are recorded with the same PR Number or JR Number. We define a *network* as the set of persons implicated in the same criminal affair. Obviously, such a network can contain a single person. This database evolves daily. The registration of the court cases is carried out by state agents. We have a database collecting data from individuals who have committed offenses from 2015 to 2021 with approximately 16454 records and each record is described by around fifty (50) variables or attributes. For ethical reasons, certain information that are likely to identify people are simply deleted (e.g. name, address).

In the field of international migration, these penal data allow the observation of the relations between several categories of actors: states, migrants and criminal groups. Migrants are identified by data relating to birth and nationality as well as by the qualification of the facts alleged or suffered; generally, these are offenses against the public peace, relating to the entry and stay of foreigners, or offenses against persons, constituting the smuggling of migrants or the trafficking of human beings. These two variables, origin and observed facts, can be crossed. The age is used for minor migrants implicated or most often victims. The gender variable is of real interest for the study of HTNs, while the profession variable reveals some occupations often related to potential trafficking.

### B. Supervised Classification Methods

Supervised Classification is one of the main techniques offered by machine learning algorithms. Typically, a supervised classification algorithm works in two phases. In the first phase, called the *learning phase*, the algorithm is provided with examples of labeled data, i.e., data associated with a class. Then, once the algorithm is trained and deemed efficient, comes the *prediction phase*, where the algorithm is provided with unlabeled or unseen examples, and tries to predict the associated class.

For our case, the question we want the classifiers to answer is: Is this network surely (*Suspicious*), probably (*Probably suspicious*) or not involved in human trafficking (*Not suspicious*)?

To solve this task, we considered six classification algorithms including: Decision Tree [22] [13], Random Forest [3], Gradient Boosting [9], Logistic Regression [15], Support

Vector Machine (SVM) [5] and K-Nearest Neighbors (KNN) [16].

## C. Explainable Artificial Intelligence

Explainable Artificial Intelligence (XAI) has emerged as a fundamental field with the aim of allowing human users to understand and trust the results produced by AI algorithms [1]. In some domains, knowing *why* a given result was predicted by the algorithm can be considered as important as getting the result itself. This is particularly the case in our HTNs detection model designed to protect human trafficking victims, especially women and children. Indeed, a prediction error can induce dramatic human consequence, such as considering a victim of human trafficking as responsible for a criminal offense.

In our proposed framework, the experts don't only want to predict whether a given network might be *Suspicious*, *Probably suspicious* or *Not suspicious*, but also want to understand why the algorithm chose to classify a given example the way it did, thus leading to a better understanding of the studied phenomenon.

In this paper, we considered the explanation of the Decision Tree and Random Forest classification models. For Decision Trees, they can be easily interpreted, with a *direct reason*, derived by traversing and displaying the unique root-leaf path covering the given network. We also provide the abductive explanations (AXps) for the Random Forest as defined in [11].

We borrow the definition of a machine learning (ML) classification model from [11]. It is defined by a set of features $\mathcal{F} = \{1, \ldots, m\}$, and by a set of classes $K = \{c_1, c_2, \ldots, c_k\}$. Each feature $j \in \mathcal{F}$ takes values from a domain $D_j$. The feature space is defined as $\mathbb{F} = \prod_{j=1}^{m} D_j$. The notation $x = (x_1, \ldots, x_m)$ refers to an arbitrary point in feature space, whereas the notation $v = (v_1, \ldots, v_m)$, with $v_i \in D_i$, $i = 1, \ldots, m$, refers to a specific point in feature space. An instance (or example) denotes a pair $(v, c)$, where $v \in \mathbb{F}$ and $c \in K$. An ML classifier is characterized by a classification function $\tau$ that maps the feature space $\mathbb{F}$ into the set of classes $K$, i.e. $\tau : \mathbb{F} \to K$. Given an ML model, a computing a classification function $\tau$ on feature space $F$, a point $v \in F$, with prediction $c = \tau(v)$, with $v = (v_1, \ldots, v_m)$. A *PI-explanation* (AXp) is any minimal subset $\mathcal{X} \subseteq \mathcal{F}$ such that:

$$\forall (x \in \mathbb{F}).[\bigwedge_{i \in \mathcal{X}} (x_i = v_i)] \to (\tau(x) = c)$$

In [11], the authors proposes a purely propositional encoding for computing explanations of the random forests classification outputs, thus enabling finding PI-explanations with a SAT solver. The details on how the AXps are generated for a given Random Forest, we refer the reader to the paper by Izza et al. [11].

Let us note that for the remaining classification methods considered in this paper, to the best of our knowledge, no possible explanations is currently available for the output.

## D. Resampling Techniques

There are two main approaches to random resampling for imbalanced classification; they are oversampling and undersampling.

- Random Oversampling: Randomly duplicate examples in the minority class. It involves randomly selecting examples from the minority class, with replacement, and adding them to the training dataset.
- Random Undersampling: Randomly delete examples in the majority class. It involves randomly selecting examples from the majority class and deleting them from the training dataset.

Both approaches can be repeated until the desired class distribution is achieved in the training dataset, such as an equal split across the classes. Figure 2 illustrated both random resampling techniques.

In this paper, we used: RandomUnderSampler, RandomOverSampler and Synthetic Minority Oversampling TEchnique for Nominal and Continuous (SMOTE-NC) [4]. Unlike Synthetic Minority Oversampling TEchnique (SMOTE), SMOTE-NC can be used for dataset containing numerical and categorical features which is the case of our reformulated dataset that we will describe in the following section. All of these techniques are implemented in a python toolbox to tackle the curse of imbalanced datasets in machine learning [12].
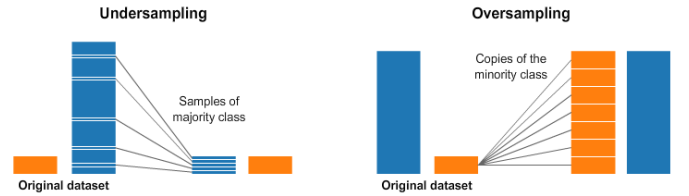


Fig. 2. Random resampling techniques

## IV. APPLICATION TO HUMAN TRAFFICKING NETWORKS DETECTION

### A. Data reformulation - Feature Engineering

Feature engineering is a crucial but labor-intensive component of machine learning applications [2]. Data scientists spend much of their effort designing preprocessing pipelines and data transformations [2]. The new features might be ratios, differences or other mathematical transformations of existing features.

In our study, it is worth noting that we have done some prior experiments between different scaling methods such as: MinMaxScaler [18], StandardScaler [18] and Log Transformation [8]. Experimental results showed that using StandardScaler achieved the highest performances on our dataset in terms of the classification algorithms used in our task. For this reason, we will use *StandardScaler* [18] as our scaling method.

StandardScaler is one of the simplest scalers which transforms numerical features to have a mean of 0 and a variance of 1. It accomplishes this by first centering (i.e. subtract the

mean from each value in the distribution) and then scaling (i.e. divide each result by the standard deviation). It can be calculated via the given formula:

$$x_{\text{scaled}} = \frac{x - \mu}{\sigma}, \quad \mu = \frac{1}{n}\sum_{i=1}^{n} x_i, \quad \sigma = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(x_i - \mu)^2}$$

where $x_{\text{scaled}}$: the scaled value of the original value $x$, $x$: the original value, $\mu$: the mean of a set of $n$ values, $\sigma$: the standard deviation of a set of $n$ values, $n$: the total number of values in the set, and $x_i$: the $i$-th value in the set of $n$ values. This method will apply on our numerical dataset on each column individually except the boolean column which contains only 0 and 1.

The original legal database is reformulated following two major steps:

1) *Extracting potential HTNs*: In this step, we aim to recognize the potential HTNs and to gain an understanding of the complexity of sex workers' situations by first querying the legal database to identify all the women. Secondly, for each woman, we look in the legal database for all the people involved in the same criminal case i.e. the same PR Number or JR Number. The output of this step is a HTNs database which contains potential HTNs and then it was annotated by an expert in human trafficking with three different classes: *"Not suspicious"*, *"Suspicious"* and *"Probably suspicious"*.

2) *Reformulating the HTNs database as a numerical HTNs database*: To determine the set of numerical features, we follow the methodology used by the experts in human trafficking to annotate the potential HTNs extracted in step one. First, they considered a so-called *suspected list* (this term will be used later in the paper) of offenses among all the offenses which related to human trafficking such as: *"not registered in the health file"*, *"pimping"*, *"lack of a health booklet"*, *"soliciting on the public road"*, *"illegal residence"*, *"corruption of a minor"*, *"illegal operation of a liquor store"* and *"exhibition of films or images contrary to morality"*. Indeed, a network is labeled as *"Suspicious"* thanks to a fine observation of the age structure, the plurality of offenses in the *suspected list* for which the individuals are implicated and the diversity of their countries of origin. While a *"Probably suspicious"* network is the one who involves a smaller number of people or only a single person potentially belonging to a case of human trafficking without the whole people being present in that case and who prosecuted for at least one of the offenses in the *suspected list* above. For this reason, *"Probably suspicious"* needs to be further verified by the judicial procedure to make the final decision. Finally, a network is labelled *"Not suspicious"* when the age structure is less diverse (e.g. only older women, only older women and men with the absence of young girls or only older

women and a young girl with the absence of a man etc.). and especially there is no one prosecuted for at least one of the offenses in the *suspected list* as well as there is no man prosecuted for pimping. From this methodology, we deduced a vector of 5 essential features for each case:

- <u>nbGirls</u>: The numerical feature which counts the presence of young girls (age between 19 to 25 years old). This feature is an indicator of the suspicious nature of the network.
- <u>nbWomen</u>: The numerical feature which counts the presence of older women (age between 30 to 35 years old). It shows that the network is well established over time.
- <u>nbAunts</u>: The numerical feature which counts the older women (aunts - age greater than 40 years). They usually play the role of guardians.
- <u>nbOffenses</u>: The numerical feature which stands for the number of different offences belonging to the *suspected list* presented above.
- <u>atLeastOneManPimp</u>: The boolean feature which indicates the presence of at least one man prosecuted for pimping is a high indicator of human trafficking. It has the value 1 if a man exists else 0 if a man does not exist.

We finally obtained a numerical HTNs database. It contains 377 court cases with 333 court cases assigned as *"Not suspicious"*, 16 as *"Suspicious"*, and 28 as *"Probably suspicious"* which is an imbalanced dataset.

### B. Experimental Evaluation

In a multi-class classification with imbalanced data (i.e. the overwhelming number of examples from the majority class (or classes) will overwhelm the number of examples in the minority class), to overcome this problem, we will use the resampling techniques described in section III-D.

*1) Comparative evaluation:* The comparative experimental evaluation between different classification methods including: Decision Tree, Random Forest, Gradient Boosting, Logistic Regression, SVM and KNN are conducted. It is worth noting that for each classifier, we have done some preliminary experiments on the numerical HTNs dataset by using GridSearch and Repeated Stratified 10-fold Cross Validation with 3 repetitions for each configuration and take the mean of Macro-Averaged F1-score as the result (see Figure 1(a)). Table I shows the best hyperparameters configuration for each classifier.

To evaluate the performances between those classifiers (see Figure 1(b)), the numerical HTNs dataset was splitted into training set (70%) and testing set (30%) using the stratified sampling to maintain the same class distribution in each subset. We repeated this operation 100 times over the dataset, each time by randomly shuffling the order of samples before splitting. Then, each time, all the classifiers with their best hyperparameters are trained on training set and evaluated on testing set. After 100 iterations, we reported the mean of empirical results by using different metrics such as: Accu-

| Models | Hyperparameters setting | Best hyperparameters |
|---|---|---|
| Decision Tree | 'max-depth': [1, 2, 3, 4, 5], 'min-samples-leaf': [1, 2, 3], 'min-samples-split': [2, 3, 4, 5] | 'max-depth': 4, 'min-samples-leaf': 2, 'min-samples-split': 2 |
| Random Forest | 'max-depth': [5, 6, 10, 20], 'min-samples-split': [2, 3, 4, 5], 'n-estimators': [100, 200, 500, 1000] | 'max-depth': 6, 'min-samples-split': 5, 'n-estimators': 200 |
| Gradient Boosting | 'max-depth': [1, 2, 3, 4, 5, 6], 'learning-rate': [0.1, 0.5, 1], 'n-estimators': [50, 100, 150, 200, 250, 300] | 'max-depth': 2, 'learning-rate': 0.5, 'n-estimators': 200 |
| Logistic Regression | 'C': [0.01, 0.1, 1.0, 10, 100], 'max-iter': [4000, 5000, 10000, 20000], 'penalty': ['l1', 'l2', 'elasticnet', 'None'], 'solver': ['newton-cg', 'lbfgs', 'sag', 'saga'] | 'C': 1.0, 'max-iter': 4000, 'penalty': 'l1', 'solver': 'saga' |
| SVM | 'kernels': ['poly', 'rbf', 'sigmoid', 'linear'], 'C': [0.001, 0.01, 0.1, 1.0, 10, 100] | 'kernels': 'rbf', 'C': 1.0 (default) |
| KNN | 'n-neighbors': [2, 3, ..., 20], 'weights': ['uniform', 'distance'], 'metric': ['eclidean', 'manhattan', 'minkowski'] | 'n-neighbors': 5, 'weights': 'distance', 'metric': 'manhattan' |

racy, Macro-Averaged Precision, Macro-Averaged Recall and Macro-Averaged F1-score.

All the experiments (include finding their best hyperparameters in Table I) have been conducted with a computer equipped with Intel(R) Core(TM) i9-10900 CPU @ 2.80GHz with 62Gib of memory. We obtained the experimental results as shown in the Table II (these results can be reproduced using a fixed seed number for the random generator). First column of Table II lists the different settings between the combination of resampling techniques, the scaling method and hyperparameters tuning. The second column shows the classification methods. The remaining columns show the values of quality metrics obtained as a result of applying different settings and classification algorithms used in the first and second column respectively.

From the analysis of the obtained results, that is shown in Table II, it can be seen that applying scaling method improved the performances for most of the classifiers across different resampling techniques, indeed, it has a slight improvement for tree-based algorithms like Decision Tree, Random Forest and Gradient Boosting because they are fairly insensitive to the scale of the features but it significantly improved the performances of SVM and KNN by the fact that all SVM kernel methods are based on distance as well as KNN which uses distance-based algorithms to calculate the distance between data points to determine their similarity, as that being said, they require the preprocessing steps for scaling the features. Using Hyperparameters Tuning also increased the performances of at least 3 classification algorithms (SVM used its default hyperparameters) for each resampling technique when comparing each of its 3rd setting with its 2nd setting. *Random Undersampling* technique which randomly delete examples from the majority class leads to poor performances for most of the classifiers compared to other resampling techniques except SVM which seems to perform well on small data with the combination of the scaling method. Other classifiers

seem to suffer from the lack of the training data as we reduced the number of samples from the majority class and thus lead to underfitting. However, results obtained from the classification algorithms performed on *Random Oversampling* are still less efficient compare to *Imbalanced data* apart from KNN. Using *SMOTE-NC* technique on imbalanced dataset outperformed other resampling techniques at least 4 out of 6 classification algorithms in terms of Macro-Averaged F1-score across different settings. The highest F1-score was achieved by applying KNN with the scaling method, its best hyperparameters configuration, and *SMOTE-NC*.

*2) Explaining classification results for understanding HTNs:* To better explain the experts of HTNs, we will use a model of Decision Tree to predict a given court case of HTNs into one of the three classes. It should be noted that using the original features or the scaling features for Decision Tree, we have found the same hyperparameters configuration.

Let us now consider a case which contains three young girls, one older woman with the age between 30 to 35 years old. This case involves in two different offenses from the *suspected list* and there exists a man prosecuted for pimping. We can represent the given case as an instance below:

*['nbGirls': 3, 'nbWomen': 1, 'nbAunts': 0, 'nbOffenses': 2, 'atLeastOneManPimp': 1]*



Fig. 3. Decision Tree of a predicted network

According to the given instance and the provided Decision Tree's structure with its best hyperparameters using the original features, we traverse the decision tree and we can see that our model predicted this network as *Suspicious* as depicted in Figure 3 by considering only two features including *"nboffenses"* and *"nbGirls"* as a *direct reason*.

For Random Forest with its best hyperparameters (Table I), we also provide the abductive explanations (AXps) generated by the RFxpl's tool in [11]. It is also worth mentioning that, to provide the intuitive explanation for the end-users, after each explanation provided by RFxpl's tool, we scaled back each feature to its original representation.

Consider the same instance as above, to predict the given instance as *Suspicious*, RFxpl provides the explanation as below:
*"IF nbOffenses = 2.0 AND atLeastOneManPimp = 1.0 THEN class=Suspicious"*.

As we can see, RFxpl predicted this case as *"Suspicious"* with the explanation that it contains two different offenses in the *suspected list* and there exists at least one man prosecuted for pimping.
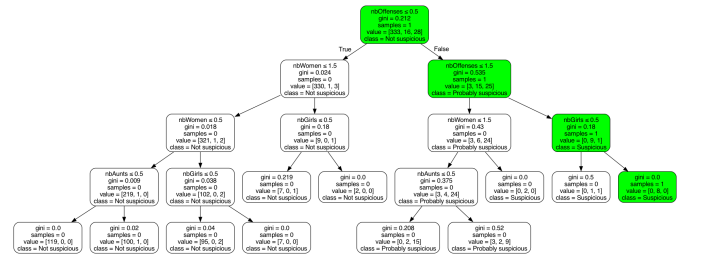
TABLE II

RESULTS OF CLASSIFICATION USING DIFFERENT SETTINGS APPLIED ON DIFFERENT CLASSIFIERS

| Different Settings | Models | Accuracy | Macro-Averaged Precision | Macro-Averaged Recall | Macro-Averaged F1-score |
|---|---|---|---|---|---|
| (1st Setting) Imbalanced Dataset Features Scaling: No Hyperparameters Tuning: No | Decision Tree | 0.9582 | 0.8443 | 0.7784 | 0.7878 |
| | Random Forest | 0.9623 | 0.8788 | 0.8025 | 0.8140 |
| | Gradient Boosting | 0.9609 | 0.8657 | 0.7895 | 0.8014 |
| | Logistic Regression | 0.9615 | 0.8755 | 0.7900 | 0.8018 |
| | SVM | 0.9595 | 0.8709 | 0.7878 | 0.8031 |
| | KNN | 0.9439 | 0.8163 | 0.6742 | 0.6920 |
| (2nd Setting) Imbalanced Dataset Features Scaling: Yes Hyperparameters Tuning: No | Decision Tree | 0.9582 | 0.8443 | 0.7790 | 0.7886 |
| | Random Forest | 0.9627 | 0.8803 | 0.8058 | 0.8173 |
| | Gradient Boosting | 0.9609 | 0.8657 | 0.7895 | 0.8014 |
| | Logistic Regression | 0.9647 | 0.8800 | 0.8113 | 0.8217 |
| | SVM | 0.9614 | 0.8772 | 0.8001 | 0.8164 |
| | KNN | 0.9531 | 0.8369 | 0.7301 | 0.7376 |
| (3rd Setting) Imbalanced Dataset Features Scaling: Yes Hyperparameters Tuning: Yes | Decision Tree | 0.9592 | 0.8425 | 0.7777 | 0.7864 |
| | Random Forest | 0.9646 | 0.8804 | 0.8149 | 0.8277 |
| | Gradient Boosting | 0.9613 | 0.8602 | 0.8004 | 0.8107 |
| | Logistic Regression | 0.9653 | 0.8828 | 0.8164 | 0.8268 |
| | SVM | 0.9614 | 0.8772 | 0.8001 | 0.8164 |
| | KNN | 0.9591 | 0.8594 | 0.7762 | 0.7864 |
| (1st Setting) Random UnderSampling Features Scaling: No Hyperparameters Tuning: No | Decision Tree | 0.9252 | 0.7813 | 0.7645 | 0.7494 |
| | Random Forest | 0.9172 | 0.8058 | 0.7985 | 0.7731 |
| | Gradient Boosting | 0.9295 | 0.8014 | 0.7996 | 0.7774 |
| | Logistic Regression | 0.9276 | 0.8164 | 0.8055 | 0.7844 |
| | SVM | 0.9368 | 0.7972 | 0.8254 | 0.7939 |
| | KNN | 0.8850 | 0.7329 | 0.6230 | 0.6406 |
| (2nd Setting) Random UnderSampling Features Scaling: Yes Hyperparameters Tuning: No | Decision Tree | 0.9253 | 0.7814 | 0.7652 | 0.7503 |
| | Random Forest | 0.9175 | 0.8057 | 0.8005 | 0.7747 |
| | Gradient Boosting | 0.9296 | 0.8018 | 0.8002 | 0.7781 |
| | Logistic Regression | 0.9304 | 0.8118 | 0.8108 | 0.7846 |
| | SVM | 0.9518 | 0.8355 | **0.8450** | 0.8274 |
| | KNN | 0.9267 | 0.8245 | 0.7694 | 0.7593 |
| (3rd Setting) Random UnderSampling Features Scaling: Yes Hyperparameters Tuning: Yes | Decision Tree | 0.9385 | 0.8048 | 0.7661 | 0.7588 |
| | Random Forest | 0.9297 | 0.8237 | 0.8095 | 0.7924 |
| | Gradient Boosting | 0.9132 | 0.7656 | 0.7911 | 0.7520 |
| | Logistic Regression | 0.9425 | 0.8351 | 0.8153 | 0.7992 |
| | SVM | 0.9518 | 0.8355 | **0.8450** | 0.8274 |
| | KNN | 0.9256 | 0.7878 | 0.7571 | 0.7460 |
| (1st Setting) Random OverSampling Features Scaling: No Hyperparameters Tuning: No | Decision Tree | 0.9426 | 0.7974 | 0.7723 | 0.7597 |
| | Random Forest | 0.9483 | 0.8186 | 0.7929 | 0.7841 |
| | Gradient Boosting | 0.9478 | 0.8120 | 0.7859 | 0.7764 |
| | Logistic Regression | 0.9570 | 0.8262 | 0.8053 | 0.8011 |
| | SVM | 0.9505 | 0.8249 | 0.8218 | 0.8075 |
| | KNN | 0.9461 | 0.8009 | 0.7930 | 0.7779 |
| (2nd Setting) Random OverSampling Features Scaling: Yes Hyperparameters Tuning: No | Decision Tree | 0.9428 | 0.7975 | 0.7736 | 0.7611 |
| | Random Forest | 0.9489 | 0.8198 | 0.7976 | 0.7877 |
| | Gradient Boosting | 0.9478 | 0.8120 | 0.7859 | 0.7764 |
| | Logistic Regression | 0.9546 | 0.8199 | 0.7974 | 0.7929 |
| | SVM | 0.9488 | 0.8168 | 0.8128 | 0.7961 |
| | KNN | 0.9518 | 0.8084 | 0.7935 | 0.7866 |
| (3rd Setting) Random OverSampling Features Scaling: Yes Hyperparameters Tuning: Yes | Decision Tree | 0.9444 | 0.8231 | 0.7807 | 0.7679 |
| | Random Forest | 0.9507 | 0.8285 | 0.8040 | 0.7948 |
| | Gradient Boosting | 0.9473 | 0.8060 | 0.7922 | 0.7797 |
| | Logistic Regression | 0.9534 | 0.8160 | 0.7961 | 0.7905 |
| | SVM | 0.9488 | 0.8168 | 0.8128 | 0.7961 |
| | KNN | 0.9522 | 0.8099 | 0.7791 | 0.7762 |
| (1st Setting) SMOTE-NC Features Scaling: No Hyperparameters Tuning: No | Decision Tree | 0.9604 | 0.8548 | 0.7910 | 0.8006 |
| | Random Forest | 0.9626 | 0.8602 | 0.8072 | 0.8139 |
| | Gradient Boosting | 0.9632 | 0.8701 | 0.8079 | 0.8180 |
| | Logistic Regression | 0.9604 | 0.8427 | 0.8174 | 0.8147 |
| | SVM | 0.9560 | 0.8270 | 0.8109 | 0.8076 |
| | KNN | 0.9595 | 0.8433 | 0.8400 | 0.8304 |
| (2nd Setting) SMOTE-NC Features Scaling: Yes Hyperparameters Tuning: No | Decision Tree | 0.9604 | 0.8545 | 0.7912 | 0.8018 |
| | Random Forest | 0.9640 | 0.8713 | 0.8186 | 0.8268 |
| | Gradient Boosting | 0.9630 | 0.8732 | 0.8102 | 0.8206 |
| | Logistic Regression | 0.9594 | 0.8441 | 0.8163 | 0.8151 |
| | SVM | 0.9566 | 0.8429 | 0.8179 | 0.8153 |
| | KNN | 0.9618 | 0.8524 | 0.8445 | 0.8386 |
| (3rd Setting) SMOTE-NC Features Scaling: Yes Hyperparameters Tuning: Yes | Decision Tree | 0.9614 | 0.8632 | 0.7947 | 0.8032 |
| | Random Forest | 0.9657 | 0.8758 | 0.8241 | 0.8330 |
| | Gradient Boosting | 0.9633 | 0.8700 | 0.8148 | 0.8250 |
| | Logistic Regression | 0.9588 | 0.8386 | 0.8168 | 0.8143 |
| | SVM | 0.9566 | 0.8429 | 0.8179 | 0.8153 |
| | KNN | **0.9666** | **0.8883** | 0.8324 | **0.8434** |

In order to be well distinguished between these three different classes, let us consider 2 more examples.

The first case contains only a single young girl prosecuted for one of the offenses in the *suspected list*. The instance below presents this case:

*['nbGirls': 1, 'nbWomen': 0, 'nbAunts': 0, 'nbOffenses': 1, 'atLeastOneManPimp': 0]*

RFxpl provides the prediction with the explanation as below:

*"IF nbGirls = 1.0 AND nbWomen = 0.0 AND nbOffenses = 1.0 AND atLeastOneManPimp = 0.0 THEN class=Probably suspicious".*

As we can see, RFxpl predicted this case as *"Probably suspicious"* with the explanation that it contains only one young girl prosecuted for one of the offenses in the *suspected list* without the whole people being present in that case (i.e. the absence of older women and/or a man in the same case). As a result of this kind of classification, this case needs further investigation or verification by judicial procedure in order to make the final decision.

The second case contains a young girl with an older woman with the age greater than 40 years old. The instance below presents this case:

*['nbGirls': 1, 'nbWomen': 0, 'nbAunts': 1, 'nbOffenses': 0, 'atLeastOneManPimp': 0]*

RFxpl provides the prediction with the explanation as below:

*"IF nbOffenses = 0.0 AND atLeastOneManPimp = 0.0 THEN class=Not suspicious".*

As we can see, RFxpl predicted this case as *"Not suspicious"* with the explanation that it contains no one prosecuted for the offenses in the *suspected list* as well as it does not exist a man prosecuted for pimping. Indeed, this case contains only a young woman and an older woman with the age greater than 40 years old. For this reason, this case is not a human trafficking network.

## V. ANALYSIS AND DISCUSSION

### A. Misclassifications made by classification models

In this section, we present the types of misclassifications made by different classifiers on our numerical HTNs database. As we can see in Figure 4 as well as the equivalent percentages of precision and recall for each class as shown in Table III, those classifiers have very similar ability of predicting the class *"Not suspicious"*. Decision Tree and Logistic Regression have the lowest percentages in term of the recall of the class *"Suspicious"* but this situation is less problematic as we will explain in the next paragraph. Gradient Boosting and KNN have the same performances, they both have at least the same or better percentages for all classes compared with the other classifiers. However, all the classifiers made the same mistakes by incorrectly classifying one *"Suspicious"* case and three of the *"Probably suspicious"* cases as *"Not suspicious"*.

In real world application of HTNs, it is less severe whether the cases of class *"Not suspicious"* and *"Probably suspicious"* are predicted in a class *"Suspicious"*, similarly, for the cases

of the class *"Not suspicious"* and *"Suspicious"* predicted in a class *"Probably suspicious"* because both *"Suspicious"* and *"Probably suspicious"* classes will further need to be verified by the experts in the domain. In contrast, it is more severe if the cases of class *"Suspicious"* or *"Probably suspicious"* are predicted in a class *"Not suspicious"*. In overall, as that being said in the previous paragraph, each one of our classifiers made 4 main classification errors (i.e. 4 False Positive of the class *"Not suspicious"*) which were considered as more severe situations.

Finally, for the application of HTNs, our models will be deployed as an aid to the prosecutor's decision in the judicial procedure by providing the reasons generated by Decision Tree and Random Forest for each of the predicted output.

TABLE III
PRECISION (P) & RECALL (R) PER CLASS CALCULATED FROM CONFUSION MATRICES. NS: NOT SUSPICIOUS, S: SUSPICIOUS AND PS: PROBABLY SUSPICIOUS

| Models | $P_{NS}$ | $R_{NS}$ | $P_S$ | $R_S$ | $P_{PS}$ | $R_{PS}$ |
|---|---|---|---|---|---|---|
| DecisionTree | 0.9880 | 0.9909 | 0.9166 | 0.6875 | 0.7741 | 0.8571 |
| RandomForest | 0.9880 | 0.9939 | 0.9285 | 0.8125 | 0.8571 | 0.8571 |
| GradientBoosting | **0.9880** | **0.9939** | **1.0** | **0.8125** | **0.8620** | **0.8928** |
| LogisticRegression | 0.9880 | 0.9909 | 0.9090 | 0.6250 | 0.7812 | 0.8928 |
| SVM | 0.9880 | 0.9939 | 0.9230 | 0.7500 | 0.8275 | 0.8571 |
| KNN | **0.9880** | **0.9939** | **1.0** | **0.8125** | **0.8620** | **0.8928** |

### B. Analysis and interpretation of results by human trafficking experts

The results provided by our tool to identify HTNs for sexual exploitation have been thoroughly analysed and discussed with researchers from HSS and experts in human trafficking. The explanations provided by our classification tool was widely appreciated and allowed a better understanding of the phenomenon. This joint analysis allowed, among other things, to sketch a first model of the structure of the potential HTNs and to understand the complexity of prostitution situations on which they are based. Figure 6 shows the birthplaces of the women and girls arrested in Saly and prosecuted in the same case for not having a health and social record or for soliciting[5]. The interweaving of migratory scales and prostitution situations helps us to understand the diversity of recruitment locations and the variety of activities and statuses according to the origin of the victims. Young girls born near Saly, in the groundnut basin, are essentially clandestine prostitutes. Among those from other regions of Senegal are women from Ziguinchor, presented as "Tantes"[6]; indeed, the victims of trafficking themselves are increasingly used by traffickers to recruit new victims; being of the same sex fosters trust. These "Madam's" aim to pay off their debt and hope to end their exploitation. At the same time, young girls born in Guinea-Bissau and Guinea-Conakry state that they have come with a

---

[5]This is a case listed in 2016 in the Register of Complaints and Minutes of the Regional Court of Thiès.

[6]Tantes are usually older women who manage the day-to-day sexual exploitation of younger girls. English-speaking girls also use the term "Madam's".
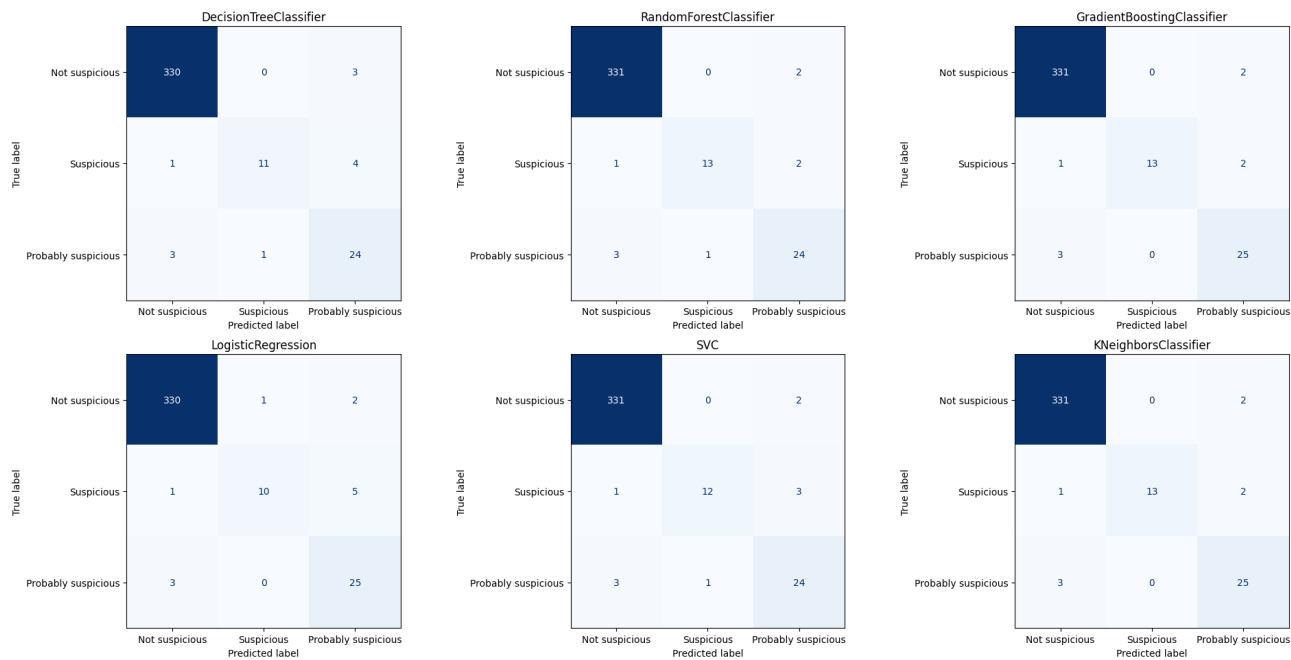
Fig. 4. Confusion Matrices present the misclassifications made by each classifier on numerical HTNs database

"guardian" and live with a "fiancé"; for them, in addition to the offence of not having a health and social security card, there is the offence of soliciting, vagrancy and illegal residence; they declare themselves to be "hairdressers"; this is a screen activity intended to hide their sexual exploitation by criminal groups.

In order to move from the specific to the general, an analysis of all the criminal cases with a combination of facts and actors similar to the case described above reveals the structure of the trafficking networks in Saly. From one case to another, the facts and actors may change but the structure of the networks remains the same. Thus, Figure 5 presents the organisational laws of a trafficking network for sexual exploitation:

- five to ten young girls, aged between 19 and 25, illegal foreigners or Senegalese; the former have a level of education equivalent to secondary school; the latter have not attended school; all declare themselves to have no profession;
- These girls are "supervised" by two older women with different profiles: one woman, aged around 30, is a hairdresser and an illegal foreigner; her nationality is that of one of the foreign girls; the second, aged between 40 and 50, is from Nigeria or Casamance, declares herself to be a prostitute and is being followed for pimping and/or corrupting minors under the age of 21;
- These "tantes" manage the young girls on behalf of the "tutors" and "fiancés"; there are one or two of them per network, aged between 30 and 35, Nigerians and/or Italians, and they are prosecuted for procuring, corruption of minors, incitement of minors to debauchery and criminal association.
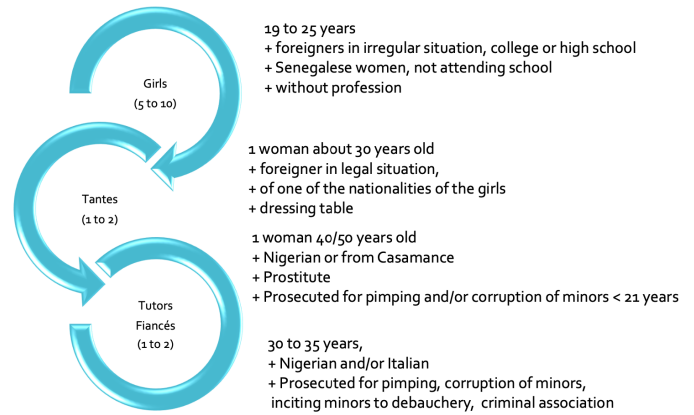


Fig. 5. Structure of trafficking networks for sexual exploitation

Thus, the trafficking network is nourished by interactions between actors, but also relies on specific territories that it highlights; this is what Figure 6 reveals to us. The places where the acts were committed, combined with the homes declared by the girls, provide a map of the activity of the trafficking networks in Saly. They are essentially spread over four neighbourhoods: Saly Carrefour, at the entrance to the city; Saly Vélingara, a precarious housing area, close to the tourist center; Saly Tapé, located in the heart of the seaside resort; and finally Saly Niakh-Niakhal, which is home to luxury residences.

The development of sex tourism in Saly explains the transit function assigned to it by the trafficking networks; the victims are not only exploited there but also "initiated" into the practices of the European prostitution market; this cannot

be the case in mining sites, conflict zones and cross-border markets. Criminal groups draw from among the young girls exploited on these peripheral sites those whom they intend to send to Saly. This complementarity of locations, according to a variable temporality, underlines the strategic position of Senegal on the trafficking routes linking sub-Saharan Africa to Europe.
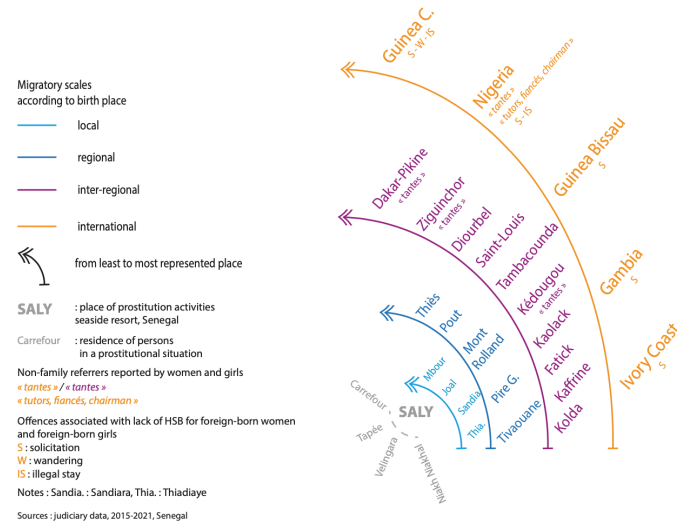


Fig. 6. Structure of trafficking networks for sexual exploitation

## VI. Conclusions and Future Work

In this paper, we make an original use of legal database to identify HTNs both sexual abuse victims and exploiters by first seeking new and relevant features expressing relationships between people involved in the same court case to form a numerical dataset. Then, six classification algorithms combined with resampling techniques alongside different settings are used to perform the comparative experimental evaluations. Our experimental results show that applying resampling technique called *SMOTE-NC* on imbalanced dataset outperformed other resampling techniques at least 4 out of 6 classification algorithms in terms of Macro-Averaged F1-score. To help the end-users, to understand the displayed HTNs, for Decision Tree and Random Forest, we also provide the explanations involving the important features that were used by those classifiers to explain how a given network was labeled.

In the future, we will seek for more relevant features from the original database for HTNs and seek more insightful explanations for the end-users. Finally, we recall that our main goal in designing a HTNs detection tool is to protect human trafficking victims, especially women and children.

The tool was presented to the UNODC in Dakar (2022), as well as to the National Unit for the Fight against Trafficking in Persons (CNLTP) of Senegal, who were very enthusiastic and confirmed the interest of the process. The software prototype was presented to CNLTP and the Senegalese Judicial Training Center (CFJ), discussed its principles and validated. Discussions with the IRD are underway to allocate seed grant to finance the development of the prototype and its deployment in Senegal.

## References

[1] Barredo Arrieta, A., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., Garcia, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R., Herrera, F.: Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. Information Fusion **58**, 82–115 (2020)

[2] Bengio, Y., Courville, A., Vincent, P.: Representation learning: A review and new perspectives. IEEE Transactions on Pattern Analysis and Machine Intelligence **35**(8), 1798–1828 (2013). https://doi.org/10.1109/TPAMI.2013.50

[3] Breiman, L.: Random forests. Machine Learning **45**(1), 5–32 (2001)

[4] Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P.: Smote: synthetic minority over-sampling technique. Journal of AI Research **16**, 321–357 (2002)

[5] Cortes, C., Vapnik, V.: Support-vector networks. Machine learning **20**(3), 273–297 (1995)

[6] Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. CoRR **abs/1810.04805** (2018)

[7] Esfahani, S.S., Cafarella, M.J., Pouyan, M.B., DeAngelo, G.J., Eneva, E., Fano, A.E.: Context-specific language modeling for human trafficking detection from online advertisements. In: Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL. pp. 1180–1184 (2019)

[8] Feng, C., Hongyue, W., Lu, N., Chen, T., He, H., Lu, Y., Tu, X.: Log-transformation and its implications for data analysis. Shanghai archives of psychiatry **26**, 105–9 (04 2014). https://doi.org/10.3969/j.issn.1002-0829.2014.02.009

[9] Friedman, J.H.: Greedy function approximation: a gradient boosting machine. Annals of statistics pp. 1189–1232 (2001)

[10] Hernández-Álvarez, M.: Detection of possible human trafficking in twitter. In: 2019 International Conference on Information Systems and Software Technologies (ICI2ST). pp. 187–191 (2019). https://doi.org/10.1109/ICI2ST.2019.00034

[11] Izza, Y., Marques-Silva, J.: On explaining random forests with SAT. In: Zhou, Z. (ed.) Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021. pp. 2584–2591 (2021)

[12] Lemaître, G., Nogueira, F., Aridas, C.K.: Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning. Journal of Machine Learning Research **18**(17), 1–5 (2017)

[13] Leo Breiman, Jerome Friedman, C.J.S.R.O.: Classification and Regression Trees. Chapman and Hall/CRC (1984)

[14] Li, L., Simek, O., Lai, A., Daggett, M.P., Dagli, C.K., Jones, C.: Detection and characterization of human trafficking networks using unsupervised scalable text template matching. 2018 IEEE International Conference on Big Data (Big Data) pp. 3111–3120 (2018)

[15] McCullagh, P., Nelder, J.A.: Generalized Linear Models. Chapman Hall / CRC, London (1989)

[16] Mucherino, A., Papajorgji, P.J., Pardalos, P.M.: k-Nearest Neighbor Classification, pp. 83–106. Springer New York, New York, NY (2009)

[17] ONUDC: The effects of the covid-19 pandemic on trafficking in persons and responses to the challenges (2021)

[18] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. Journal of Machine Learning Research **12**, 2825–2830 (2011)

[19] Tong, E., Zadeh, A., Jones, C., Morency, L.: Combating human trafficking with multimodal deep models. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL. pp. 1547–1556 (2017)

[20] UNODC: Global report on trafficking in persons (2020)

[21] Wang, L., Laber, E.B., Saanchi, Y., Caltagirone, S.: Sex trafficking detection with ordinal regression neural networks. CoRR **abs/1908.05434** (2019)

[22] Wu, X., Kumar, V., Quinlan, J.R., Ghosh, J., Yang, Q., Motoda, H., McLachlan, G.J., Ng, A., Liu, B., Philip, S.Y., et al.: Top 10 algorithms in data mining. Knowledge and information systems **14**(1), 1–37 (2008)