

University Degree in Audiovisual Systems
Academic Year (e.g. 2014-2019)

Bachelor Thesis

“Violent event detection from acoustic signals”

Óscar Otero Martínez

Carmen Peláez Moreno
Madrid, February 20, 2020



[Include this code in case you want your Bachelor Thesis published in Open Access University Repository]

This work is licensed under Creative Commons **Attribution – Non Commercial – Non Derivatives**

SUMMARY

Nowadays, the availability of big data resources has allowed the development and implementation of different algorithms to learn patterns and behaviours. Useful information can be extracted from data in order to find solutions to real world problems. This branch of study is known as machine learning. Within this field, an area of research that has largely evolved during the last years is the one related to image and audio learning. Given the multimedia resources that are currently available, this type of data has become one of the main sources of knowledge with a potential to be transformed in successful results for several situations. Some examples for the works belonging to this field are speech recognition, Music Information Retrieval (MIR), Multimedia Event Detection (MED), Acoustic Scene Classification (ASC), Acoustic Event Detection and Classification (AED/C).

One of the tasks that has been considered is the Acoustic Violent Detection (AVD). An example of a common application for this research line is related to the detection of violence in videos or movies in order to categorize them appropriately for parental control. Also, several works have been done related to video surveillance topic. Our purpose is to work in a new branch related to violence against women, more specifically, in gender-based violence centered in Intimate Partner Violence (IPV). For this reason, the *UC3M4Safety Team* has been created as a joint effort from a handful of research groups in a project called EMPATIA¹. Its aim is to develop a protocol to implement machine learning methods to prevent this type of actions.

As a contribution to the project, the task of classifying audio events trying to distinguish between violent and non-violent has been addressed in this work. In order to collect the necessary amount and type of data, we decided to work with *AudioSet* database, which comprises around 2 million annotated videos and a total of 527 classes. Also, we wanted to define a new concept for gender-based violence applied to a certain victims' situation by developing a system in which the user can select the type of categories that are more common in her daily violent situations. So the classes considered in the classification task can be personalized.

In order to characterize the data selected, we attempted to find a robust feature extraction method that goes beyond low-level descriptors (LLD), that allow us to find a feature space in which violent events were easier to distinguish. With this purpose, we have used a novel type of feature that can be obtained by running deep learning algorithms with the desired data, also known as *embeddings*. These are extracted by making use of the

¹*protEciónn integral de las víctimas de violencia de género Mediante comPutaciónn AfecTiva multi-modAl*, funded by Department of Research and Innovation of Madrid Regional Government (Y2018/TCS-5046)

transfer learning concept with the CNN model named as *VGGish*, designed by Google and provided with the Audio Set database.

For the classification task, we have taken advantage of different algorithms as SVM, CNN and LSTM. We have tried them in a multiclass classification with single label data and then use them in a preprocessing stage for a posterior binary classification between violent and non-violent events performed with a simple SVM classifier.

Keywords: **Violent Event Detection, feature embeddings, multiclass and binary classification, Neural Networks**

DEDICATION

Firstly, I would like to thank my tutor *Carmen Peláez Moreno* for all her help and commitment throughout this work. I hope the EMPATIA project achieves all its goals and is applied to real life situations in the foreseeable future.

A big thank you to all the friends I have made during my university years. Specially, I would like to point out *Irina Antich Moreno* and *Daniel Barrejón Moreno* for being such good classmates, friends and always being helpful to solve any kind of doubt.

To the two people that I will always be grateful, *my parents*. Thank you for making this possible, allowing me to study in this University and always being supportive under any circumstances.

To my older sister *Lucía*, who has always been a role model since I was kid. No matter the pass of the years, everything will always be like in the beginning.

To my girlfriend *Ana*, who has been there for me during all these years, caring so much about me everyday and always making me happy.

Finally, to my grandparents, specially *Antonio* and *Leonor*, who have always worked very hard to make their families prosper. My sister and I are the final result of those efforts.

CONTENTS

1. INTRODUCTION	1
1.1. Context	1
1.2. Objectives	1
1.3. Regulatory framework	2
1.4. Socio-economic environment	3
2. STATE-OF-THE-ART	4
2.1. Acoustic Scene Classification and Acoustic Event Detection and Classification	4
2.2. Features and methods for ASC and AED/C	4
2.2.1. Data augmentation	10
2.2.2. Transfer learning	11
2.3. Violent Event Detection	12
2.3.1. Gender-based violence	13
2.3.2. Our point of view	14
2.4. Databases	15
3. RESOURCES AND METHODS FOR ACOUSTIC VIOLENT DETECTION . .	17
3.1. Database: AudioSet	17
3.1.1. What is AudioSet	17
3.1.2. Ontology	17
3.1.3. Labels quality	20
3.1.4. Data access	20
3.2. Feature extractor	21
3.2.1. Artificial Neural Networks (ANN) and Convolutional Neural Networks (CNN)	21
3.2.2. Visual Geometry Group (VGG) model	27
3.2.3. VGGish model	29
3.3. Methodology	31
3.3.1. Synthetic Minority Over-sampling Technique (SMOTE)	31
3.3.2. Support Vector Machine (SVM)	32

3.3.3. Recurrent Neural Network and Long Short Term Memory	33
4. OUR APPROACH FOR ACOUSTIC VIOLENT DETECTION	38
4.1. Why VGGish	38
4.2. Our approach	39
4.2.1. Input data	39
4.2.2. Extracting embeddings	41
4.2.3. Assessing the differences between the two types of data access	42
4.3. Our system	47
5. EXPERIMENTS.	50
5.1. Input data preparation	50
5.2. Implementations and results	53
5.2.1. Implementation 1: SVM classifier for a multiclass classification	54
5.2.2. Implementation 2: SVM multiclass classification and SVM binary classification	57
5.2.3. Implementation 3: LSTM for multiclass classification	59
5.2.4. Implementation 4: LSTM multiclass + SVM for a final binary classification	62
5.2.5. Implementation 5: CNN for multiclass classification	63
5.2.6. Implementation 6: CNN multiclass + SVM for a final binary classification	66
5.3. Comparison and results	67
6. CONCLUSION	72
7. FUTURE WORK	74
8. PLANNING AND ECONOMIC RESOURCES	76
8.1. Planning and time organization	76
8.2. Economic resources	77
BIBLIOGRAPHY.	81
A. APPENDIX METRICS	
B. APPENDIX K-FOLD CROSS VALIDATION	
C. APPENDIX CATEGORICAL CROSS-ENTROPY	
D. APPENDIX SPECTROGRAM AND MEL SCALE	

LIST OF FIGURES

2.1	Framework for ASC addressing the idea of BOF and using a GMM to obtain statistical representation of the low-level features. $S_{\Lambda q}$ are the original labelled samples in the training set and $S_{n,\Lambda q}$ are multimodal distributions for these samples. $x_{n,\Lambda q}$ are the features extracted with an operator T which are taken by an operator S to learn the global model M . Then, in the testing process, a likelihood measure G is used in order to classify the <i>new</i> sample.	7
2.2	Difference between traditional machine learning (a) process and feature learning (b) [47]	12
3.1	First two layers of Audio Set Ontology [65]	19
3.2	Comparison between shallow networks and deep neural networks [74] . .	22
3.3	Movement of the kernel represented as a red block along the width and the height of the original input image. It is clear that it finally occupies the whole depth of the input image. The arrow indicates an approximation of the movement that the filter follows [76].	24
3.4	Max-pooling operation with a filter size of 2 and a stride of 2 for a single depth filter. Each colour represents the action portion for each operation [75]	26
3.5	Representation of a typical CNN model [80]	27
3.6	VGGish architecture	29
3.7	SMOTE algorithm [90]	32
3.8	Visualization of a hyperplane set by SVM and other decision boundaries [92]	33
3.9	Scheme or a recurrent neural network about how the loop works	33
3.10	Example of long-term dependencies situation. If the predicted output h_{n-1} depends on x_0 and x_1 , the vanishing gradient problem may appear.	34
3.11	LSTM module with the representation of the different operations that take place inside of it	35
3.12	Forget gate	35
3.13	Input gate	35
3.14	Cell state	36

3.16 Output gate	37
4.1 Flowchart about selecting violent classes	40
4.2 Confusion matrices	43
4.3 Architecture to see how the different embeddings work	44
4.4 t-SNE results for wav format with a legend that shows the labels of the corresponding samples in the original 128D space	46
4.5 t-SNE results for <i>.tfrecord</i> format with a legend that shows the labels of the corresponding samples in the original 128D space	47
4.6 Block diagram of the whole model including both types of classification. We do not extract the embeddings, but this is a diagram of the whole model considering that someone did extract them.	49
5.1 Bar plot that shows the number of samples for each of the selected classes. Clearly, the violent categories are much less populated than the others, that did not have to much any semantic criteria	52
5.2 Number of observations for train, validation and test subsets used in the different experiments	52
5.3 Confusion matrices for SVM multiclass classification for train and validation sets	54
5.4 Confusion matrices for SVM multiclass classification for train and validation sets	55
5.5 Confusion matrix for the test set for the SVM multiclass classification .	56
5.6 Accuracy values for the three sets for the SVM multiclass classification .	56
5.7 Confusion matrix for the train and validation sets for the SVM multiclass + SVM binary classification	57
5.8 Confusion matrix for the test set for the SVM multiclass + SVM binary classification	58
5.9 Accuracy values for the three sets for the SVM multiclass + SVM binary classification	58
5.10 LSTM architecture	59
5.11 Confusion matrices for LSTM multiclass classification for train and validation sets	60
5.12 Confusion matrix for the test set for the LSTM multiclass classification .	61
5.13 Accuracy values for the three sets for the LSTM multiclass classification .	61

5.14	Average and standard deviation for the train and validation set for the LSTM multiclass + SVM binary classification	62
5.15	Confusion matrix for the test set for the LSTM multiclass + SVM binary classification	63
5.16	Accuracy values for the three sets for the LSTM multiclass + SVM binary classification	63
5.17	Confusion matrices for CNN multiclass classification for train and validation sets	64
5.18	Confusion matrix for the test set for the CNN multiclass classification . .	65
5.19	Accuracy values for the three sets for the CNN multiclass classification .	65
5.20	Average and standard deviation for the train and validation set for the CNN multiclass + SVM binary classification	66
5.21	Confusion matrix for the test set for the CNN multiclass + SVM binary classification	67
5.22	Accuracy values for the three sets for the CNN multiclass + SVM binary classification	67
5.23	Test confusion matrices	70
A.1	Example of confusion matrix	
A.2	Confusion matrix for a multiclass classification [112]	
B.1	K-fold cross-validation scheme. The data is first split into train and test, and then the train is split again with this method [115]	
C.1	Cross-entropy loss function when the true label is equals to 1. As the probability of the predicted class approaches to 1, the loss function tends to zero. However, if it the probability is closer to 0.0, then the loss function increases heavily [116].	
D.1	Audio signal in time domain	
D.2	STFT [120]	
D.3	Spectrogram of an audio signal. The x-axis corresponds to time, the y-axis to frequency and the colour of the plot to the amplitudes by following the colour bar on the right.	

LIST OF TABLES

2.1	Examples of time and spectral audio features.	6
2.2	Table of studied databases	16
3.1	Fields per category in the ontology.	19
3.2	VGG ConvNet configurations.	28
4.1	Chosen classes for a small classification. <i>Screaming, Crying, sobbing</i> and <i>Gunshot, gunfire</i> are considered as the violent ones.	42
4.2	Accuracy values for audio and <i>.tfrecord</i> files	44
5.1	Relation of <i>violent</i> and <i>non-violent</i> selected classes	51
5.2	Train, validation and test set for different experiments	53
5.3	Adaptation of the results for SVM multiclass classifier to binary	57
5.4	Adaptation of the results for LSTM multiclass classifier to binary	62
5.5	Adaptation of the results for CNN multiclass classifier to binary	66
5.6	Accuracy results for the three different algorithms and the three sets for the multiclass approach	68
5.7	Accuracy results for the three different algorithms and the three sets for the binary approach	70
8.1	Relation of estimated time and actual time invested per task in hours . . .	76
8.2	Relation of estimated time and actual time invested per task in hours . . .	77

1. INTRODUCTION

1.1. Context

Violence against women has been considered the most common way of violating human rights in the whole world. This type of actions are the origin of injury and a risk factor for many psychological and physical health problems. Also, it increases the possibility of suffering from other types of health issues as physical disability, abuse of consumption of drugs and alcohol and depression [1][2]. This type of violence is usually based on deep social and cultural roots and it is undoubtedly linked to unbalanced relationships between men and women in different situations and contexts, such as economics, politics and religion [3]. Also, victims are often women who endured violence during their childhood and felt socially isolated.

One of the most common types of violence against women is the Intimate Partner Violence (IPV) which can be defined as physical and emotional damage and sexual abuse by an intimate partner with a controlling attitude. Although violence between partners in a heterosexual relationship can be exerted from both sides and also may appear in homosexual cases, the most usual situation of IPV is from men to women [4].

A way to retrieve information of what is happening in a certain situation is by making use of audio resources. Within this field, plenty of techniques have been developed in order to extract usable information for solving real world problems. As will be explained in detail in 2.3, one of these tasks consists in detecting violent situations for several applications.

In order to include the violence against women as a new branch in the violent detection, the *UC3M4Safety Team* from Universidad Carlos III de Madrid has presented the project EMPATIA² in which, among other types of data, acoustic information extracted from the victims' environment is used to perform this task [5].

1.2. Objectives

At the beginning of this project, the main objective of the work was to design and implement a system that allowed to detect and identify violent events from acoustic signals. In order to define a way of how to address the problem, one of the idea was to work in a sound detector. However, how do we know exactly what kind of sounds do we want to detect? Then, we chose the path of defining a system based on a classifier of acoustic events

²protEcciónn integral de las víctimas de violencia de género Mediante comPutaciónn AfecTIva multimedAl funded by Department of Research and Innovation of Madrid Regional Authority (reference Y2018/TCS-5046)

that gives as final output a binary tag identifying the input act as violent or non-violent to improve on the definition of violence. This is basically the global objective of the whole project. The following list collects the different partial objectives considered necessary to obtain the final one.

- To find a data resource, for example, a public database, that offers a sufficiently high number of data observations, especially, that belong to a certain type of classes that can be considered as violent in different situations.
- To discover a set of features or a feature space in which it is feasible to obtain a clearer separation between violent and non-violent events than the one obtained with conventional feature extraction methods.
- To find a customizable definition of violence that fully adapts to the victim's perception of a violent situation.
- To implement a system that allows the previous customization.
- To find an algorithm that allows to take advantage of the new features in order to perform a multiclass classification to distinguish between violent and non-violent events.
- To adapt the output obtained from the multiclass task of the model to a binary classification in order to finally identify violent actions.

1.3. Regulatory framework

When paying attention to the final application of the whole EMPATIA project, the regulatory framework ranges from aspects related to private data protection to the integral law of gender-based violence, specifically, the conditions related to the rights of the Victims' rights as well as the aggressors'.

Since the objective of the whole project consists in collecting information from day-to-day life of the victim in order to control and record her emotional state and check for signs of violence inside her usual environment, this is a continuous access to her privacy. So, all types of actions that take place in every kind of situations could be recorded. Also, when considering evidences in trials, a good point would be to present as a valid evidence data collected with this system. This is actually a good important application of the collected data due to, in several occasions, the judicial cases become stack because of the lack of credible evidences [5]. In the case of achieving a resolution for the case, the data could also help to distinguish between the different types of sexual crimes that can be carried out [6], and so the imposed penalties for the aggressor which is different for each case.

The scope of this topic within the project is so complex that a significant part of the UC3M4Safety team is specifically working towards correctly dealing with the regulatory framework. However, the scope of the task related to this specific project is still far from this type of considerations. Since the database used in this work is publicly available as YouTube videos, the legal requirements we deal with are nothing but the privacy policy of this website [7]. It may have a repercussion on the development tasks when trying to download the multimedia data, since there are some videos that are not available maybe because of this reason.

1.4. Socio-economic environment

The implementation of a system like this has its most important repercussion on the increase of women empowerment and gender equality. There have been a lot of improvements in this field mainly during the last years but still the discrimination is present in all the countries in the world.

As a part of the Sustainable Development Goals (SDG) as a program by the United Nations (UN), goal 5 is the one related to the objective of *achieving gender equality and empower all women and girls* [8]. There are several abuses that are still present and should be eliminated in order to aim to live in a more prosperous, sustainable and peaceful world. For example, 1 out of 5 women have declared to have suffered physical or sexual damage by an intimate partner. Also, 49 countries in the contemporary world do not count with a legislation to protect women from domestic violence.

It is mainly in this atmosphere of domestic violence in which the developed task can find a more important application within the whole project by defining a new protocol for the protection of the victims of gender-based violence including all the involved agents in order to implement a technological solution [5].

2. STATE-OF-THE-ART

Since the main topic of this work, acoustic violent scene detection, has not been researched as such in the literature before, we decided to start our search by addressing close two topics that have been largely studied for the last years related to acoustics scenes and events also including some aspects on violence detection mainly aimed at the detection of violent scenes in movies and video-surveillance.

2.1. Acoustic Scene Classification and Acoustic Event Detection and Classification

Acoustic Scene Classification (ASC) refers to the association of an audio sequence to a certain semantic label that describes the environment in which it took place [9]. With this idea in mind, the classification of acoustic scenarios has been tackled with two different kinds of concepts: soundscape cognition, i.e. understanding how the human being perceives the sounds subjectively from the physical environment that surrounds them [10], and Computational Acoustic Scenes Analysis (CASA), that is working on new computational methods that may help automatize this task through machine learning and processing signal techniques [11]. This notion can have many applications, such as content recognition –by allowing devices to obtain benefits and information from its situation [12]–, for medical usage [13], as a tool to aid musical recognition [14] or as a complement to Computer Vision (CV).

Simultaneously to these advances in the ASC, another related area has evolved during the last years: Acoustic Event Detection and Classification (AED/C). It can be described as the processing or treatment of sound signals in order to convert them into significant descriptions that match a listener’s sensing of the events and sources composing the acoustic environment [15]. The detection part consists in identifying the events in a temporal stream of audio and labelling them. The result is usually accompanied by the time interval in which the occurrence is set. However, classification is a task that acts directly on the event that has been already isolated and has the purpose of designating a label or class to the sound [16]. These techniques also have plenty of applications, e.g., in the medical field [17], in biological topics such as bird noise detection [18], and for multimedia information retrieval from video sources in social media [19], etc.

2.2. Features and methods for ASC and AED/C

As in many fields, the boundaries between features and methods used for audio tasks are becoming blurred due to the recent rise of deep learning based methods. Classically, a problem is addressed with a pre-processing stage of the data and then the model or method is implemented. Right now, these two stages are sometimes maintained but also

have been mixed or changed depending on how the algorithm works.

In every machine learning or pattern recognition task, for the system to be able to infer and extract conclusions from the input given, a pre-processing stage is necessary to make some transformations to the data so that they could be understood by the model. This stage is known as feature extraction and the goal is to convert the original information into a set of values or vectors that characterize the data regarding some desired properties [20].

There are several ways that allow us to perform this processing stage. The most common and basic one consists on extracting features that are closely related to the original signal which are called low-level descriptors (LLD) [21]. These are computed by performing some mathematical operations or formulas to the original data and can be considered rudimentary when comparing with other techniques. However, they are really extended and still in use nowadays [22].

In the audio field, there are two types in which all LLD can be grouped into. One of them is for the features that have been computed by considering the audio signal in its original form in the recording, i.e., in time-domain, that is the reason why they are known as time-domain audio features. The other case refers to those characteristics that are obtained from the signal after having been transformed into the frequency domain. These are commonly known as frequency-domain or spectral audio features. For the procedure of feature extraction, the signal is usually divided into frames that can be overlapped by using a sliding window, so the calculations are done per frame, obtaining a final matrix with size *number of frames* \times *number of features* [20]. It must be taken into account that the goal of the whole system is going to be fundamental at the time of deciding which features must be computed. For example, not the same features are to be used for speech recognition than for musical information retrieval. In table 2.1, some of the most popular time and spectral features are included.

Moreover, either time-domain or spectral features can also be split in two groups depending on the way they were extracted. There are the ones called short-term features. This way of computing follows the framing process explained above. The signal is initially divided in various frames that are commonly overlapped and for each frame the feature is calculated [20]. The other case refers to long-term features. This way the different computations to obtain the descriptors are performed along longer segments, sometimes the whole audio is taken [25]. For a certain long sequence, the computation of the short term features are calculated and then some statistics are performed in order to characterize this information for the long term. Actually, the long-term sequence can be defined as a set of statistics that resumes the short-term features for a sequence. The portion of audio selected for the long-term usually presents a uniform behaviour within the total data [20].

To sum up, a typical procedure would first divide the signal into frames. The next step would be to compute the short-term features considering some time-domain or spectral characteristics within the audio data that are thought to be determinant in the required

Feature	Description
Time	
Energy Entropy	It is useful to detect sudden changes from the energy of a signal. To calculate this value for a certain subframe, it is necessary to first compute the normalized energy of the subframe with respect to all the frames energy [23].
Short time energy	It is the energy for a short segment of signal. It is normally used in speech-related tasks in order to identify voiced form non-voiced fragments [24], among other things.
Zero Crossing Rate (ZCR)	This can be defined as the number of times the amplitude of the signal crosses the zero line per unit of time, i.e., changes from negative to positive. It is sometimes computed by the number of zero-crossings by the amount of samples in the frame [23].
Frequency	
Spectral Flux (SF)	It is computed to measure the spectral changes between two successive frames. To do so, the difference of their spectral energy obtained from the FFT is computed [24].
Spectral Rolloff	This represents the skewness of the shape of the spectrum given the frequency below which a concrete percentage of the magnitude distribution of the frequency transform is concentrated [24].
Spectral Centroid (SC)	This is a measure related to the spectral position. It is defined as the center of gravity of the spectrum. [23].
MFCC	It is a feature that it is widely used because it gives good results in many tasks as for example, speech recognition since it interprets the frequency bands in a very similar way to human perception while doing a separation of the fine structure of the spectra (that corresponding with the harmonics) and the coarse (i.e. the filter representing the vocal tract). It is computed from the STFT [24]. A wider explanation can be found in appendix D.

Table 2.1. EXAMPLES OF TIME AND SPECTRAL AUDIO FEATURES.

task. Finally, a series of mathematical and statistical transformations are performed in order to obtain a LLD vector per sample in order to pass these input data to a statistical model that will be trained with the objective of summarizing the properties of the signals

and develop a classification rule so as to be able to assign a faithful category to a new unlabelled observation [26]. For example, some of the long-term features that can be obtained from the short-term ones included in 2.1 could be the variance and mean in the case of short time energy, or the relation between the maximum and the mean along the frame for the SF [23].

A common approach for the utilization of this type of LLD in the long-term way is the one known as Bag-of-frames (BOF). The term has been acquired from the text treatment field as an analogy to the technique known as "bag-of-words", which consists in treating the text data as a global and unique distribution of words without giving importance to their order inside the sentences [27]. The analogous system for audio tasks consists in representing signals as statistics of long-term sequences of their spectral local features. This method has been applied in different fields in the ASC, for example, in order to analyze "soundscapes", analogy for audio landscape, and polyphonic music [28]. It was established that this approach was a proper option in order to work with soundscapes since the order of the audio sequence inside of them was not really relevant for its understanding. However, in the music task, the ordering of the parts of the sound was determinant to retrieve the musical information, so the technique was discarded for this topic [29]. A common way of exploiting this concept of BOF is by using the already mentioned MFCC and the really well-known technique in the literature as GMM. These are generative methods in which the feature vectors are treated as they were produced by a multimodal distribution which consists on a sum of Gaussian distributions [9]. In figure 2.1, it is shown the structure of a Bag-of-frames approach.

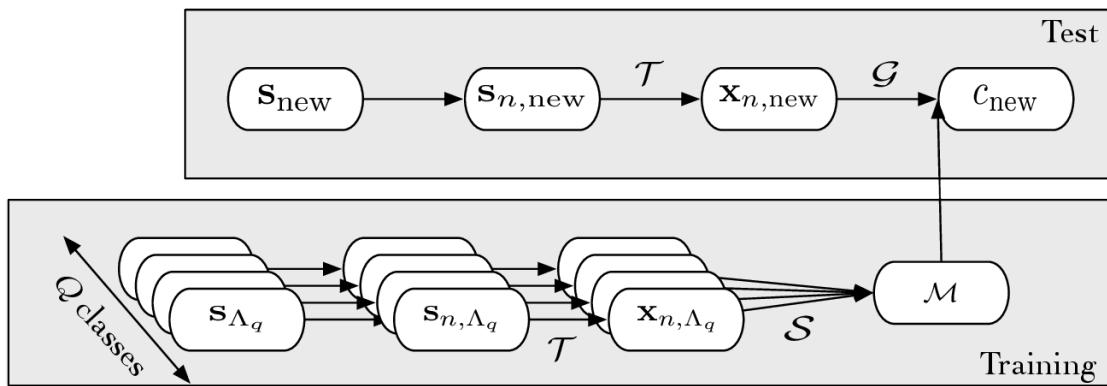


Fig. 2.1. Framework for ASC addressing the idea of BOF and using a GMM to obtain statistical representation of the low-level features. S_{Λ_q} are the original labelled samples in the training set and S_{n,Λ_q} are multimodal distributions for these samples. x_{n,Λ_q} are the features extracted with an operator T which are taken by an operator S to learn the global model M . Then, in the testing process, a likelihood measure G is used in order to classify the *new* sample.

A similar way of acting can be found for AED. In this case, the typical procedure has been found to be defined as a two-steps process. For the case of Multimedia Event Detection (MED) [19], in the first part, some features are generated at a clip level, usually

by concatenating frame-level features, so then, in the second stage, a binary or multiclass classification can be done in order to confirm the presence of an event in a whole clip of video. This idea follows the same concept of BOF explained above since the generation of features do not attend to their order in the whole clip.

Another similar approach has been found for this task, but in this case it differs from the BOF and it is applied to real-life recordings. For the first stage, a classification of already isolated events is performed so a vocabulary of acoustic actions can be built. For this purpose, it is needed a set of short-term video recordings in which the semantic meaning of the corresponding audio event must be clear and highlighted over all the others sounds that take place along the sequence. In the second stage, the detection procedure takes place along the whole clip by making use of the previous configured vocabulary of acoustic data [30]. The way the events are modeled in this task is by using Hidden Markov Models (HMM). This is a common technique that has been put in practice plenty of times in the audio world and consists in modeling the temporal succession of the events within a longer sequence. This can be done by saving the order of the events in a transition matrix that contains the probability of one sound occurring after another one [9].

Also, this process can be combined with a classification of the events detected during the long sequence in order to make an ASC based on what type of sounds are taken place within the scene [9]. This is actually the initial idea desired for our work. Building a violent acoustic dictionary with a classification of monolabel short time recordings so as to, then, perform a detection process for violent acoustic events in order to categorize the scene and check if there is violence on it. However, we just mainly made an approach for the very first part about classifying what we considered violent events. This will be explained further on in chapters 4 and 5.

In some cases, the mentioned low-level descriptors, sometimes referred as manually-crafted features, are not enough to model data since they cannot achieve a meaningful representation of the original information for the system to learn from. Also, as mentioned above, the selection process of this type of features plays a significant role in the final output so that the criteria of the designer could be considered a limitation for the classification task [31]. For this reason, a new strategy of addressing feature extraction has appeared based on the idea of finding more specialized properties by using specific engineering algorithms that explores the given data in order to find non-human recognizable patterns. The techniques that allow to perform this task are usually called automated feature generation algorithms or machine-learnt features [32].

In order to obtain these high-level features, there is not just one standardized machine model that allow to compute them all in a particular form. The calculation method, nevertheless, differs from one approach to another. One example could be the combination of low-level data with high-level semantic descriptors that consists on the inference of diverse dimensions by using Support Vector Machines so as to find similarities among music genres. Particularly, they consider the output probabilities of the SVM classifier as

a high-level feature space in which the distance between samples from different classes can be measured [33]. This algorithm has been one of the most popular ways to track classification problems.

There are plenty of works in the literature that use this technique for audio classification tasks [34] [35] [9]. It is originally known to be a binary classifier but, nowadays, there have been some software implementations that allow its use in multiclass problems. This algorithm makes the classification by taking into account pairs of observations that are more likely to be misclassified and draws a hyperplane as the optimal boundary between the two classes [36]. Since it has been used for the experiments in this work, a more detailed explanation will be included in section 3.3.

Several methods that have been really successful on the high-level feature extraction, apart from other tasks as detection and classification, lay under the umbrella of Artificial Neural Network (ANN). This concept was inspired by human biology and how the neurons communicate among them in the brain in order to interpret the input or sensory data humans collect [37].

In the field of AED these have been used in order to implement techniques also for automatic learning of features. For example, a technique of boosting is used to extract discriminative features which resulted in a better performance than MFCC for real-life acoustic data [38]. They first made an attempt based on modeling the sequence data with HMM explained before and combined them with a trained ANN. Also, an approach based on SVM and GMM was tried.

The development of the ANN in different ways has made possible the implementation and creation of new algorithms during the last years that exploits the concept of neural models from different perspectives. Is the case of Convolutional Neural Network (CNN), which follows the way of working of the neural networks but performing convolution operations over the input data [36]. This has achieved really good results for image tasks in the Computer Vision field, but also implemented successfully in audio problems. As in [39], where is proposed a method based on CNN in order to classify ten acoustic scenes. In order to feed the model they used log-mel spectrogram images. More information about this type of data can be found in appendix D. A detailed explanation of CNN can be found in subsection 3.2.1.

Also, another type of networks that has been treated largely with the purpose of working with sequential data is the Recurrent Neural Network (RNN). The concept behind these models is its capability of being able to learn information from previous elements inside a sequence. An adaptive model was designed called LSTM, which does basically the same but for longer sequences. One possible approach for this algorithms could be the one presented in [19]. They first classify frame level audio against a set of semantic units. Then, the output is taken as a representation with variable length in a clip level which feed the LSTM as input. The results overtake SVM and FC in similar tasks. More information about this type of networks is included in detail in subsection 3.3.3.

Related to the development of these more complex models as CNN or LSTM, a novel type of feature has been used in the last years that differs from the already explained low-level and high-level ones: Deep Neural Network (DNN) have grown increasingly for plenty of classification tasks and so in the multimedia area. The problem with this type of systems is the huge amount of data that is needed to make them work properly, which can be translated in a lack of labelled data. One of the habits that has been currently resorted by the researchers consists of learning what is called deep data *embeddings* from extensive collections of, in our case, audio and use them so as to perform shallow classifications by using simpler datasets. There have been implemented some models about this topic, such as Look, Listen and Learn (L^3) [40] net that uses as input for the embedding extractor the linear-frequency log-magnitude spectrogram of 60 million audio samples, the system called SoundNet [41] that has been designed to obtain embeddings from training a deep audio classifier in order to predict the output of a deep image classifier and the VGGish network, designed by Google researchers. This last case is the one we used in this work and it will be explained in more detail in section 3.2.

2.2.1. Data augmentation

Deep Neural Network models, as the ones just mentioned above, have proved to achieve very good results for audio classification tasks by extracting information from audio-spectral type of features. Sometimes, the need of a huge quantity of data becomes a main problem due to the lack of enough correctly labeled samples, which result in an impossible actual employment of these deep and complex methods [42]. In order to obtain to solve this problem, some methods have been developed inside the field of data augmentation. This consists in a strategy which has been commonly adopted so as to increase the amount of data in the training process of the models and avoid some problems as overfitting³ so the model can be more robust [44].

The way to obtain more training samples can vary depending on the task that is been developed. In the case of audio data, these can subjected to some low-level deformations in order to generate observations that maintain the semantic of the original ones but are actually different for the model. In [45], they perform a couple of transformations such as time stretching of the signal, i.e. by changing the speed of the sample; pitch shifting, by raising or lowering the pitch; adding background noise and compressing the dynamic range of the audio sample. Also, other types of augmentation can be found in the literature that work in a more synthetic manner. For example, in [46], it is provided a technique based on sample mixed data augmentation for a domestic audio tagging application. Among the different proposed methods, one of them is the *SamplePairing*, which generates a synthetic new observation by computing the average between two input vectors.

³Overfitting is a term referring to the disability of the model to generalize a pattern from the training data which means that the NN achieves very different results for train and test sets [43].

In our case, we opted for a technique also based on synthetic samples creation which is called Synthetic Minority Over-sampling Technique (SMOTE). This selects a random sample from the less populated classes within the training data, pick one of its closest neighbours and generates a new observation placed between them in the feature space. It is explained in much more detailed in 3.3.1.

2.2.2. Transfer learning

Another option to work with deep models, as the ones mentioned in 2.2, consists in pre-training them with available huge collections of data so as to use them in a pre-processing stage and then apply them to models considerably much less complicated to address other type of problems. This is possible due to a kind of techniques commonly known as transfer learning and domain adaptation.

The typical consideration in plenty of machine learning tasks consists on extracting the training and testing subsets of data from the same feature space and same distribution. When one of these initial assumptions change, it is necessary to rebuild the whole model from the initial point, including new training data, which means a lot of computational cost and loss of efficiency [47]. Its working manner can be explained as an analogy of how humans transforms their ability on a certain task to obtain knowledge for other purpose. An example could be how musicians apply their previous experience to get to know faster how to play another instrument.

The first work in which this topic has been treated widely was in 1995 in the workshop *Learning to Learn* [48] and since then many approaches have arisen and renamed the same idea as knowledge consolidation or inductive transfer. However, it was 10 ten years later, in 2005, when the first idea of the ability of a system to identify and apply previously learned skills to completely new problems appeared from the hand of the Broad Agency Announcement (BAA) 05-29 of Defense Advanced Research Projects Agency (DARPA)'s Information Processing Technology Office (IPTO) [47]. This can be expressed as a relation between a *source* task, where the abilities are learned, and *target* task, the novel problem that needs to be resolved. As a difference with other similar methods, in this concept of transfer learning the roles of these two are not equal since the weight of the target is much heavier. In figure 2.2, the difference between a common machine learning approach and the use of transfer learning is included.

The same idea can be understood from a mathematical point of view as the analysis of the relationship between the two different spaces from types of targets [47].

Considering a *domain* D that is composed by a feature space denoted by X and a marginal probability named $P(X)$, where $X = \{x_1, \dots, x_n\} \in \chi$. The whole domain can be expressed as $D = \{\chi, P(X)\}$ where x_i is a certain vector inside the feature space.

In the same way, a *task* can be defined as T formed by a label space γ and an objective predictive function η . The task formulation is $T = \{\gamma, \eta\}$. This predictive function cannot

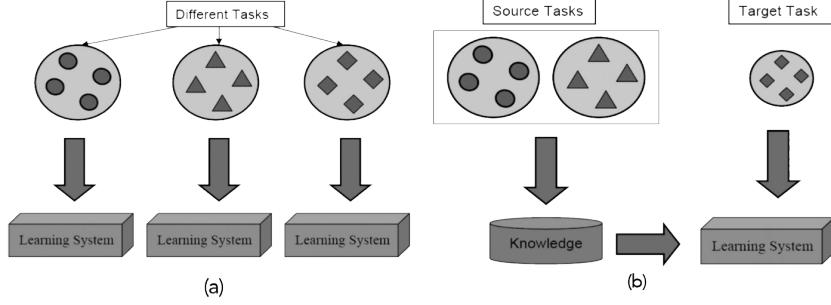


Fig. 2.2. Difference between traditional machine learning (a) process and feature learning (b) [47]

be observed, however the intention is to learn it from the training data, that is composed by pairs of the form $\{x_i, y_i\}$, $x_i \in X$ and $y_i \in \gamma$.

The predictive function η can be used to predict a corresponding label of a new sample x . From a probabilistic perspective, this new label can be expressed as $P(y|x)$. So, the task T can be defined as $T = \{y, P(Y|X)\}$, in which $Y = \{y_1, \dots, y_n\} \in \gamma$. For each vector x_i , the function η finds a prediction y_i .

Once these parameters have been defined, considering the source domain D_S , task of source domain T_S , the target domain as D_T and its respective task as T_T , the transfer learning has the purpose of obtain the condition distribution in the target domain $P(Y_T|X_T)$ with the information extracted from D_S and T_S where $D_S \neq D_T$ or $T_S \neq T_T$ [49].

We have took advantage of this technique in order to obtain features for our final model. As it is explained further on in 3.2.3, the network VGGish has been pretrained with a huge dataset with similar data as the one we take. Thanks to this method of transferring learning from models, we can extract our features without need of training the complex VGGish, which eases the whole process in a considerable way.

2.3. Violent Event Detection

All the multimedia information available can be applied to many fields and with different connotations. One of the variants that has appeared in the acoustic scenes and events field is the one applied to violence detection. For this case, an essential point before addressing any problem is to decide what kind of definition the word 'violence' is going to adopt since it is a subjective concept. An objective perspective has been given by the World Health Organization as "The intentional use of physical force or power, threatened or actual, against oneself, another person, or against a group or community, that either results in or has a high likelihood of resulting in injury, death, psychological harm, maldevelopment or deprivation" [50]. There are other definitions found in different works as "physical violence or accident resulting in human injury or pain" [51] or "any situation or action that may cause physical or mental harm to one or more persons" [23].

Recent studies have treated this problem in different ways due to all types of con-

ditions that it may take place in. During the last years, the possibility of creating and providing audiovisual content has grown widely, which has led to an enormous variety of topics in which, some of them, could be considered inappropriate for certain audiences. This is the reason why there have just been done works related to the field of video content analysis and detection of violence.

In some cases, audio and image features have been combined to address these problem [52]. However, it has been found that sound information could be really useful and a more efficient way of working compared to image, since it is easier to process and the cost is lower. Related works have utilized audio features in the time-domain and in the frequency-domain, similar to the ones explained for ASC, then combined with a normal SVM classifier [23]. Other researches have tried more complicated models with the intention of improving the classification task. It is the case of using DNN, fed with both image and audio data, which performs the task more efficiently [53]. Violence detection has also been used for other applications such as video surveillance. For example, one of the scenarios for this purpose consists on preventing violent acts inside elevators [54]. For this case, the considered dangerous situations are composed of anti-social actions that are likely to happen in this kind of places, concretely, urinating, vandalism and attacks on vulnerable victims, such as women, children or elderly. The framework proposed is based on audio-visual data, but the master classifier is driven by audio, due to the possible subtleness of the scenes that are desired to detect. So, first the audio incident detector triggers the process when a non-silent event takes place. Then, the image processing begins in order to extract information related to who is involved in the action and how aggressive is it. Another utilization of the surveillance approach is its use for the evolution of smart cities [24]. For this goal, since the system will be implemented in real-life environments, one of the advantages about working with data coming from sounds is the respect for privacy, that, otherwise, using video recordings would be violated.

The difference in these two applications, apart from the task they are addressing, resides on the data they are working with. For violent content analysis, the data usually comes from fictional audio sources as movies or video-games. However, for real-environment systems, the data is extracted straight from actual day-to-day life situations. In this second case, some disadvantages can be appreciated. For example, the signals are not pre-processed, which means the original properties of the sound are not modified so the processing part before classification becomes tougher. Also, the presence of background noise is more common and loudness of some events, as speech, may vary with time [55].

2.3.1. Gender-based violence

In the same way, in the recent times, late studies have shown that 35% of women from all over the world have been victims of physical or sexual damage [56], and 43% of women from Europe have declared going through some psychological or mental violence at least

once in their lives [57]. In this context, it is necessary to define the concept of gender-based violence, which can be described as the set of harmful behaviours that are focused on women and girls just because of their sex, such as female children and wife abuse, sexual assault, dowry-related murder and marital rape, among others.

Particularly, violence against women involves any act of verbal or physical force, extortion or lethal denial which has a woman or girl as a target and provokes the physical or psychological hurt, humiliation or irrational privation of liberty and contributes to continue women subordination [1]. Within this definition, it can be considered that most of the times that these violent situations take place, they are originated due to persons that are supposed to be part of the victims' closest circle of trust, i.e., their husbands or boyfriends. This is called Intimate Partner Violence (IPV) and it is recognized as a public health problem affecting women across their life span resulting in different undesirable unhealthy outcomes, such as depression, chronic pain and even dead [58].

2.3.2. Our point of view

As a contribution to the EMPATIA⁴ project developed by Universidad Carlos III de Madrid, the main goal in this work is to make progress in detecting gender-based violence situations, specifically applied to day-to-day scenes, in which IPV is likely to be present. One of the parts from the proposed system within the whole project is composed by wearable devices that the victim can carry to collect diverse types of information and process them to obtain conclusions and increase the efficiency. Among these accessories, we can find a pendant that senses the user's voice and the surrounding audio to analyse what is happening at a given moment. For our purpose, the interesting part resides on achieving auditory data so as to detect violent incidents that consists of sounds already known for characterizing these episodes considered dangerous by the victim.

The definition that is assigned to violence is really important in order to define which audio events should be taken into account. However, considering the subjectiveness of this concept, categorizing violence for every type of user is an extremely difficult task. For this reason, the final idea to answer this question is to make the victim able to decide which kind of hearing events the system must be aware of. In the complete project, this can be carried out by a phone user interface which displays a list of sound events and she has the labour of picking up those that are violent according to her criteria. Since the development of this tool is out of the scope of this work, we have decided to implement a simpler mechanism which will be explained in subsection 4.2.1.

⁴*protEcciónn integral de las víctimas de violencia de género Mediante comPutaciónn AfecTiva multi-modAl* funded by Department of Research and Innovation of Madrid Regional Authority, in the EMPATIA-CM research project(reference Y2018/TCS-5046)

2.4. Databases

A main objective was to find a database that allows to build a system with the desired characteristics. This should satisfy the need of a rich variety of acoustic events with an essential big representation of violent sounds. In the table 2.2 is represented a relation of the different databases that have been considered for the realization of this work.

The last three options are the ones that better suit the problem of this project. *VSD Benchmark* was the first option we chose. Within the two ways of working given, the movies and the YouTube videos, the former was the easiest to use since the annotation specified exactly what kind of violent events were present in the scene and the onset and offset time stamps within the whole film. However, this data is copyright restricted and it was necessary to pay for the contents. The latter was accessible but just indicated the presence of violence, without determining the type of event. Another choice was the *Freesound dataset* because it contains all types of videos so we could extract those classes that are more interesting for us. However, it is still under an annotation process and it is not ready to download yet. As a final conclusion, we decided to go for *AudioSet*, which will be explained further on in 3.1.

Name	Description	Considerations
URBAN-SED [42]	10,000 soundscapes with sound events. Every soundscape contains 1 to 9 sound events with strong annotations.	Events are completely specified but it just contains three interesting types of classes.
UPC-TALP [59]	It belongs to CHIL project, for the AED task. Isolated acoustic events that occur in a meeting room environment.	Payment is needed to achieve the data and the classes are a little out of our topic.
MIVIA: Audio Events Data Set for Surveillance Applications [60]	6,000 events with background noise.	The classes included belong to our topic, but they are just three: glass breaking, gun shots and screams.
TUT rare sound events [61]	Source files for creating mixtures of rare sound events (classes baby cry, gun shot, glass break) with background audio.	Similar problem to MIVIA: just from three interesting classes.
IEEE AASP Challenge [62]	Composed by ASC and AED. It is formed by two subtasks: OL (Office-live) and OS (Office Synthetic)	Labels for both subtasks are out of our scope since they are likely to happen in an office environment: keyboard clicks, hitting table, etc.
TUT-SED Synthetic 2016 [63]	Isolated sound event samples were selected from commercial sound effects	The variety of classes is large enough but for our purpose just four of them are useful.
VSD benchmark [64]	Violent events from 32 Hollywood movies and 86 YouTube web videos, together with high-level audio and video concepts.	Payment is needed to purchase the movies and the videos do not specify the type of violent event
AudioSet [65]	An ontology of 632 audio event classes and a collection of 2,084,320 human labeled 10-seconds sound clips from YouTube videos.	Our final choice. Plenty of the videos have more than one audio label but we were able to adapt the data to the problem because of the huge amount of clips.
Freesound dataset (FSD) [66]	Filling AudioSet ontology with 297,144 audio samples from Freesound.	This may seem a very good option as well but it is not available as of today.

Table 2.2. TABLE OF STUDIED DATABASES

3. RESOURCES AND METHODS FOR ACOUSTIC VIOLENT DETECTION

In this chapter we are going to explain the different resources we have chosen to create an approach that allowed us to implement our final models. First, we include a detailed description of the database that we have chosen in section 3.1. Then, we comment and analyze the system we are going to employ so as to obtain high level features for our classification experiments in section 3.2. For this point, we first make an extensive explanation of the ANN and CNN in 3.2.1, then we can explain the model VGG in 3.2.2, just to finally end in 3.2.3 with the explanation of the model we used in: VGGish.

3.1. Database: AudioSet

3.1.1. What is AudioSet

Audio Set can be described as a large-scale sound dataset with the goal of putting the availability of audio and image data on the same level. It is composed by a huge variety of 2,084,320 manually-annotated audio events in video format and it is organized by following an ontology formed by 632 different audio classes. The data has been extracted from YouTube videos and the labelling process has been based on diverse factors such as metadata, context and content analysis. It has been developed by Google with the purpose of producing an audio event recognizer that can be applied to plenty of acoustic situations coming from the real world [65].

Due to some limitations found during the building process that will be explained along the section, the current final release of the dataset is composed by 1,789,621 segments, where each video usually last 10 seconds which results in a total duration of 4,971 hours of video and audio. After executing the selection process in which the different labels were populated with the final corresponding segments, a total of 527 classes were gathered, out of which 485 counted with at least 100 samples [65] [67].

3.1.2. Ontology

In order to put this dataset together, the events have been organized in an abstract hierarchy. This is composed by higher-level classes which describe a certain type of sound event and also act as parents of other labels that refer to more specific events. With this purpose, the relationship among different classes needs to be non-exclusive, so labelling similar audio events may result into a more general class, the parent, if there is ambiguity. This is also helpful for labellers due to group the clips in an easier and faster way.

The Audio Set Ontology was collected considering some fundamental guidelines explained below:

- A complete collection of all desirable labels was prepared so that it can be used to define sound events from real-world aural data.
- When labelling an audio event the result must match the criteria of a common listener.
- Different categories should be easy to distinguish by an ordinary listener. In the case that two different labels do not satisfy this requirement, these should be merged. With this condition, the plethora of possible labels remains limited.
- The distinction of two different classes must be done by relying just on the audio, it cannot be accompanied by image or visual information.
- The hierarchy should not be very deep by keeping the number of children per parent class to no more than 10. This also eases the annotation labour.

It is easy that an ontology of this volume gets leaned or biased towards a particular direction due to several factors, such as the subjectiveness of its creators or the selection of the initial set of classes used when starting to label. With the intention of generating a primary list that covers a wide range of audio events in an objective way, the researchers decided to apply an impartial, web-scale text analysis from the very beginning. They agreed on detecting hyponyms of the word "sound" by utilizing a modified version of the famous technique called *Hearst patterns* [68], a method proposed to automatically acquire hyponymy lexical relations from unrestricted text. As a result, an enormous collection of terms came up. This was filtered by considering how well these terms represented audio, i.e. by combining together the global frequency of occurrence and how exclusively these are recognized as hyponyms of "sound" instead of other terms. As an output of the final process, a list of 3000 terms was obtained.

With respect to the hierarchical relation among categories, it was constructed by the authors with the main intention that this satisfied their human comprehension of the sounds. Event though this is a subjective manner, it also makes sense since this is how the hierarchy best performs its labour on helping human labelling. After all the organization process, the model is not based on a strict organization as a singular node can appear in many different locations, i.e., a single node can be child of different parent nodes. As we said before, the final result is composed by 632 audio event labels and, in the hierarchy, the deepest class is six level away from the top parent. Figure 3.1 shows the nodes that belong to the two top layers of the ontology.

This whole structure is offered to any user as a file in JSON format. A couple of fields have been included for each label in order to describe its meaning and make clear its position within the hierarchy. A description for all of these can be found in table 3.1.

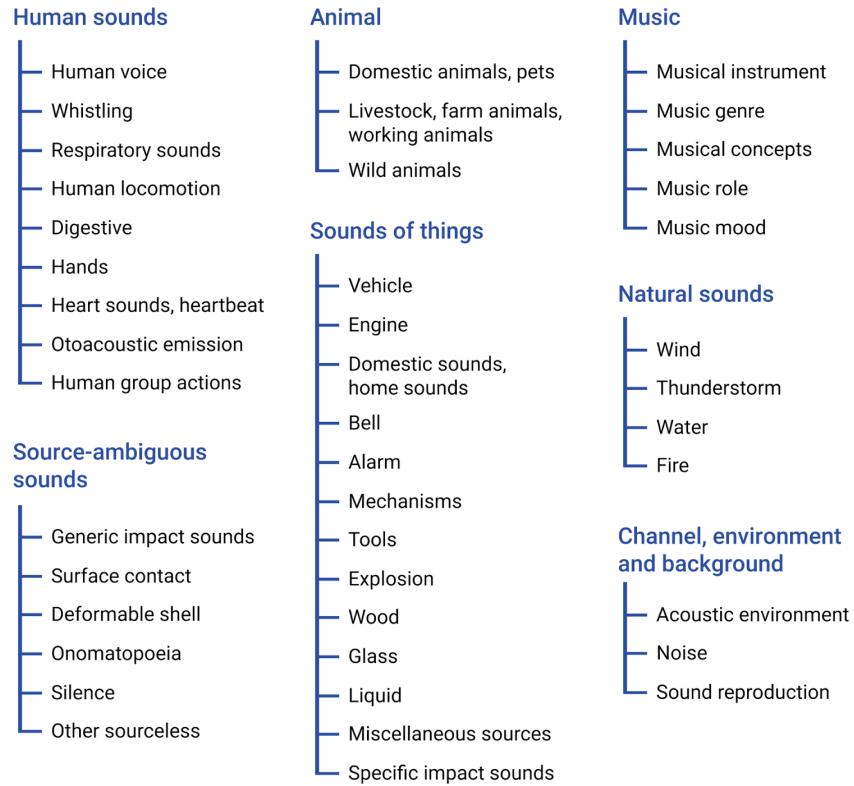


Fig. 3.1. First two layers of Audio Set Ontology [65]

ID	This field includes the Knowledge Graph Machine ID (MID) that best describes the sound or its source. It is used as a primary identifier for the class.
Display name	Short name formed by one or two words that identifies the audio class. It sometimes includes a small alternative separated by a comma so it does not feel ambiguous.
Description	One or two explaining sentences so the meaning of the category is more defined. These can be extracted from Wikipedia or WordNet.
Examples	At least one example of the label is provided as a URL of a YouTube video.
Children	An array filled by the Machine ID (MID)s from all immediate children of the class.
Restrictions	It specifies if the category in question either has been discarded or there are no audio clips under it.

Table 3.1. FIELDS PER CATEGORY IN THE ONTOLOGY.

For the field *ID*, the identifiers are known as Machine ID (MID) and belong to the *Knowledge Graph* designed by Google [69]. This is a knowledge base that Google ser-

vices use to improve the quality of its search results and it is composed with information extracted from a wide variety of sources. The MIDs are the identifiers of the different elements that belong to this huge dictionary. For instance, the MID of the word "Speech" is "/m/09x0r". Another field that deserves a special explanation is the one corresponding to *Restrictions*. Within all the categories of sound events, there are two flags that indicate an exclusive behaviour that differ from a typical label: "blacklisted" and "abstract". The former refers to a class that has been hidden from labellers due to its confusing meaning. The latter has been used for those classes that are just employed as intermediate nodes in order to provide a better grouping inside the organization, and are not expected to be used in the implementation tasks. In total, out of the 632 categories, 56 have been categorized as "blacklisted" and 22, as "abstract".

3.1.3. Labels quality

For the rating process, the video segments were proposed to the rating workers without any type of extra meta data. For each video, they had the purpose of indicating if it belonged to one or more labels by setting whether a classes was *present*, *not present* or *unsure*. The number of raters for each segment was three, and a majority-vote method was used in order to decided if the label was finally present in the video segment [65].

By the end of the process, 76.2% of the votes were unanimous and the *unsure* vote had been barely used. A small number of errors during this stage were found due to confusing labels, human mistakes and understanding silent and really soft audio sounds as different. Also, some labels were commonly misunderstood so they were removed from the ontology.

With the intention of improving the quality of the labels, quality assessment has been performed internally in order to proceed to a rerating process. Apart from giving better instructions to the raters, a procedure to tag the segments with clusters is also executed [67]. However, this stage is not completed yet, so we decided to ignore this fact and continue as the labels were equally qualified even thought this may have an effect in the final result.

3.1.4. Data access

The data can be obtained through the website [67] in two different formats:

- Files in .csv format that include for each video segment its YouTube video ID, start time, end time and the one or more labels it belongs to.
- Instead of the audio files themselves, they provide already extracted audio features for each segment in compressed files that are available in TensorFlow (TF) [70] record files that can be easily downloaded.

Together with these two options of obtaining the data, a neural network model is also provided in the database website which can be used to extract embeddings as high-level features. It is called *VGGish* [71] and it has been pre trained on a preliminary version of the database YouTube-8M [72]. A more detailed explanation of this system is later included in 3.2.3.

Based on the first choice, we did some attempts of downloading a couple of videos from YouTube using the information given in the .csv files. Then, we decided to try to extract our own embeddings features by running the VGGish network with the downloaded videos in order to see how they look and check their behaviour on a classification task. Also, for the second option, we downloaded the provided embeddings already extracted by the Google developers and see if they were also useful for our problem. A little experiment explained in 4.2.3 were implemented so as to check if the two manners of obtaining the embeddings were equivalent.

3.2. Feature extractor

VGGish model design was based in an already existed network called VGG. In order to describe the embedding features extractor we used, it is first needed to make a clear explanation of what a Artificial Neural Network is and how Convolutional Neural Network behave. Then, an explanation about how the VGG model is implemented included. Then, we can move to analyze the differences and changes included in VGGish.

3.2.1. Artificial Neural Networks (ANN) and Convolutional Neural Networks (CNN)

Before going deep in the description of CNN, the core idea of ANN must be explained. As it was already mentioned at the end of subsection 2.2, these are based on the communication way that the neurons follow inside a brain to interpret the stimuli they receive. This is actually the hidden concept behind the ANN intention, which is transforming input data into a desired output so that it can be used to perform a given task.

A basic architecture is usually composed by an *input layer* of neurons, a certain number of *hidden layers* and an *output layer*. The model receives some input, typically a feature vector, and pass it trough the hidden layers in order to perform some transformations. The neurons of a layer are completely connected to the neurons of the next layer but neurons in the same layer are not usually connected. For this case, we can say that they are Fully-Connected (FC) layers. Finally, the output layer gives the resulting outcome of the network [37]. This is just the basic and initial nature of the idea but it was implemented in many different ways depending on the input data and the desired task.

Within the different possibilities of designing an artificial neural network, there are two main options. The selection of one of them depends on the complexity of the problem, the format of the input data and how it is combined with other techniques. One of them

covers all the models that are structured with a shallow architecture. An example of this would be a basic type of network known as multi-layer perceptron (MLP) that follows the structure explained for ANN with Fully-Connected layers but composed just by one hidden layer. This type of systems have performed well in many tasks but they have shown some limitations in more complicated applications [73].

When working in real-world cases such those that involve natural signals, such as human voice, natural sounds and visual scenes, for example, it is needed the development of a deeper architecture in order to extract and understand the complexity inside this type of data so as to build. For example, visual and speech recognition tasks are exploited with system that follows the hierarchical transformations that the humans perform when understanding this type of data. reliable and useful internal representations. This type of models are collected inside the Deep Neural Network (DNN) field. DNN has been one of the most studied fields since it was proposed in 2006 [73]. In figure 3.2, it can be found a comparison between shallow and deep models.

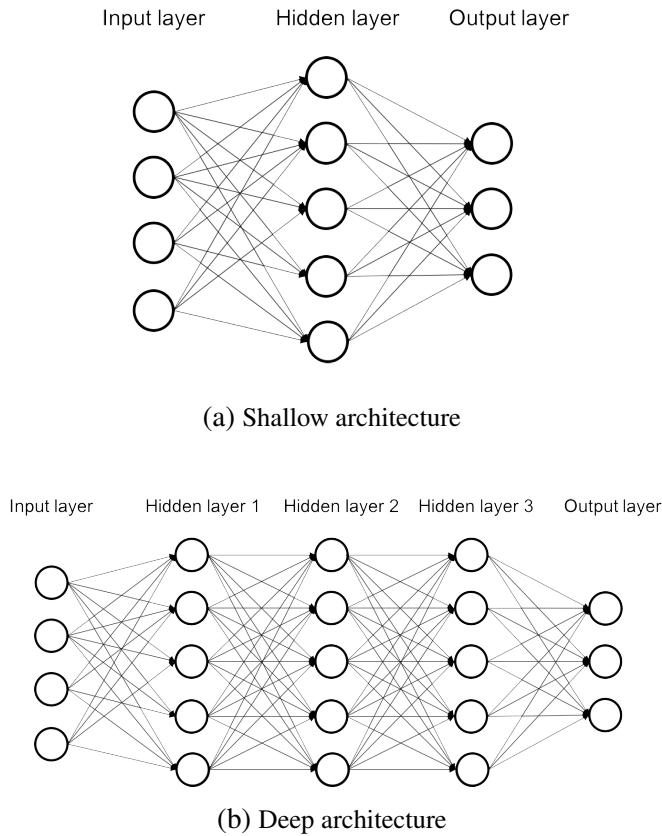


Fig. 3.2. Comparison between shallow networks and deep neural networks [74]

One of the most explored implementations inside the DNN is the one called Convolutional Neural Network (CNN). This maintain the basic idea of receiving some inputs and apply a dot product operation to them before subjecting the data to a non-linear function but with the difference that input data are expected to be images [75]. The architecture of a Convolutional Network is based on the organization of the Visual Cortex in the human brain. This means that single neurons responds to some stimuli only in a certain region of

the visual field, the Receptive Field. Then, a collection of these fields overlap with each other so as to cover the full visual space [76].

About the input data, it is introduced as a matrix instead of as a row vector. Then, the architecture of the layers adapt to the fact that input data are images. They do not consist of a single row of neurons, but they have a shape constituted by three dimensions: height, width and depth [75]. In this way, the net can capture the spatial and temporal relations present in an image by applying some filter functions. With this construction, the model can be trained with a smaller number of parameters and also reuse the weights.

A Convolutional Network can be described as a sequence of layers in which every layer converts a volume of activation results from the previous layer to another by computing a differentiable function. The architecture of a CNN is always composed by three type of distinguishable layers which are *Convolutional*, *Pooling* and *Fully-Connected*. A normal architecture will follow the next succession of layers [75]:

- INPUT: this represents the image at the entrance of the whole model. It commonly is a 3D matrix , one dimension per channel of colour, and contains the values of the original pixels. In the case of a RGB image the number of channels should be three, for example.
- CONVOLUTIONAL: considering the neurons that are connected to a region of the input image, this layer computes their output by performing a dot product between the weights and the value of the pixels. The depth of the resulting data will coincide with the number of filters used.
- POOL: this performs a fixed operation that consists of a downsampling of the spatial dimension. It typically reduces by half the height and the width of the input image.
- FC: this is the final layer which gives the output of the whole net. It is a vector-shape layer and is completely connected to the previous layer. In a classification task the number of neurons matches the amount of classes of the problem.

Also, another type of layer can be optionally added in order to improve the performance of the model. This is called the normalization layer. It is a good practice in some cases in order to put a limit for the layers output. When using an unbounded activation function, in order to bound the result of the output layers, the normalization process can be used before the activation operation is applied. Two of the most common techniques are Local Response Normalisation (LRN) and batch normalization [77]. However, in many cases the contribution of these layers has been almost null for the performance of the model [75].

Note that this list represents the order of the layers in the architecture, however, the number of layers of each type vary depending on how the model is designed and how deep it is desired to be. Some examples of deep architectures can be found later in section 3.2.2, in table 3.2.

With this building, a Convolutional Network passes the information of the pixels layer by layer while reducing the size of the input images in a way that is much easier to process. This is performed without missing any feature since all of them are essential to compute a good prediction outcome.

In the convolutional layer, a convolution operation is performed for each small region of the input image. The element that carries out this operation is called *kernel* or *filter*. This partial region that the kernel covers in the input data is set by deciding the height and the width as a designing parameter. About the depth, it must match the depth of the input data, for an image this can be understood as the number of channels, so the filter can slide through the image along its other two dimensions. The results are obtained by computing the convolution between the pixels inside the region and the values of the filter. An example can be observed in figure 3.3. For each filter, a 2D activation map is obtained after subjected the result of the convolution to a non-linear operation. One of the most common functions used for activation is the Rectified Linear Unit (ReLU) function. Finally, we will pack all these in a unique volume that will consist in the output of the convolutional layer [75]. For example, if we have a number of 64 filters, we will get then a group of 64 2D activation maps that will be the input of the next layer. Traditionally, it was the first convolutional layer the one that was responsible of capturing low-level features that can be related to edges or colours in the image, for example. The upper layers were the ones in charge of extracting high-level features so the network was able to understand the image similar to how the humans do [76].

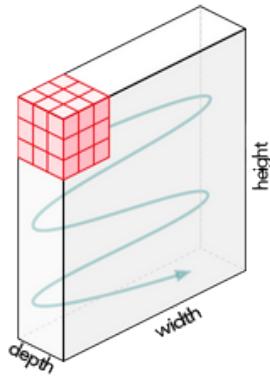


Fig. 3.3. Movement of the kernel represented as a red block along the width and the height of the original input image. It is clear that it finally occupies the whole depth of the input image. The arrow indicates an approximation of the movement that the filter follows [76].

We have included how the input interacts with the convolutional layer and also introduced the concept of depth of the resulting volume. However, it is needed to explain how the total size of this output is computed and what it depends on. There are three main parameters that play a role in this decision: *depth*, *stride* and the *zero-padding* [75].

- DEPTH: it is a hyperparameter that corresponds to the amount of kernels used in

the convolutional layer when learning from the input. The set of neurons that extent through the depth dimension operating in the same region is called depth column. This is the previous mentioned red block shown in figure 3.3.

- STRIDE: this indicates the displacement of the sliding filter along the image. If it has a value of 1 (non-strided), this means that the kernel will shift a position of one pixel along the width. When it is set to 2, then it moves two pixels. It rarely has values greater than 2. The moving process consists on the filter going from right to left with a hop determined by the stride value. When the right margin is reached, then it jumps to the beginning of the next row from the left margin and repeats the process in the same way until the image is completed [76]. The bigger the stride, the smaller the size of the output volume.
- ZERO-PADDING: this is also called just *padding*. It consists on adding zeros to the borders of the image. What is a hyperparameter in design is the size of this padding. This allows to control the width and height of the output volume so we can keep it the same through the different layers. When it is indicated a "valid padding", it means that no padding is added and the output size will be reduced. If we indicate "same padding", then it has a value of 1 and means that the resulting volume will have the size of the input [78].

As shown below, the total size of the obtained volume, D , can be computed as function of the input size, W , the spatial 2D dimensions of the convolutional filter, F , the stride of the filters shifting, S , and the amount of padding in the zero-padding, P [75].

$$D = \frac{W - F + 2P}{S} + 1$$

The convolutional layers can be grouped all together, one after another, but they are usually interpolated by what it has been previously defined as *pooling* layer. This habit has the objective of reducing the width and the height of the resulting volumes in a progressive manner in order to decrease the total number of training parameters and control the computational cost, and, therefore to avoid overfitting [75]. There are two types of pooling operation: the max-pooling and the average-pooling. The former just returns the maximum value from the portion of the image where the filter is placed. However, the latter computes the average of all the values in this portion [76]. The pooling layer acts in an independent way on every input depth level.

Two hyperparameters are needed in order to configure the spatial portion in which the pooling is computed: the filter size, related to length of height and width, and the stride. Also, for this operation zero parameters are introduced since it consists on a fixed operation over the input data. In figure 3.4, it is shown how the max-pooling operation is performed [75].

Together a convolutional layer and a pooling layer constitute a typical complete layer structure in a Convolutional Network. The number of these complete layers can vary

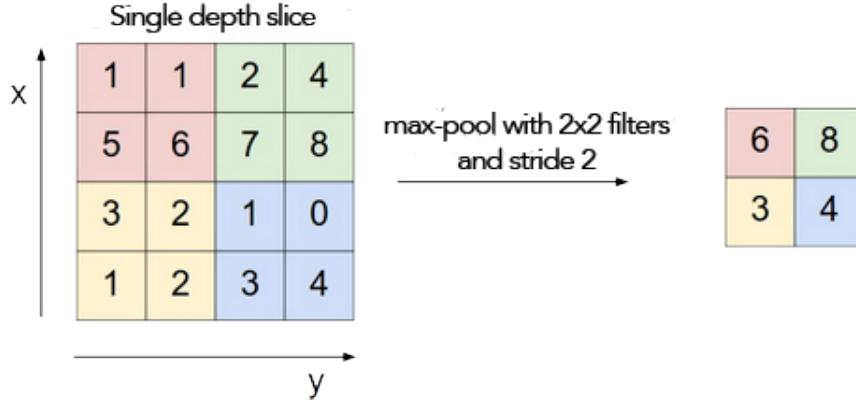


Fig. 3.4. Max-pooling operation with a filter size of 2 and a stride of 2 for a single depth filter.

Each colour represents the action portion for each operation [75]

depending on the design and the task needs. These represent the main core of the network before passing to the last layers and calculating the final output [76].

For the final step, a common way of acting is to apply some non-linear activations in order to learn from the high-level features resulted from the output volume of the last convolutional layer. This takes place in the Fully-Connected layers ending the model. This part has the shape of a Fully-Connected layer, as the one described in the basic model of ANN. To make the weights from the convolutional layers available to feed this last part of the model, it is necessary to flatten the output volume and obtain a column vector. The learning process is possible due to a backpropagation operation performed in every epoch of training. As a total output, the values that represent the classification of the different samples into one or another category are obtained [76]. A common classification technique is the one called *Softmax*. This normalizes the result of the previous FC layer into a vector whose values follow a distribution that total sum results in 1. This type of output is called *soft predictions* and represent the probability for each sample of belonging to a category in the classification [79]. However, there exists other activation functions that can compute the classification output in other ways.

In figure 3.5 an scheme of a common Convolutional Neural Network is shown. This takes an input image, compute the features along three groups of convolutional layers plus max-pooling and, finally, includes three Fully-Connected layers, being the last one a softmax function layer with as many neurons as classes that gives the soft predictions for this image of belonging to each class.

Convolutional Neural Network (CNN) were originally designed to work with images but they have been applied to other fields as audio. An image can be defined as a matrix of values, i.e. pixels, in two or three dimensions, so the only prerequisite to use a CNN is to have the input data in this form. For audio researching, it has been commonly used the log-scaled mel spectrogram, that is a way of representing audio in a scale different than frequency that is more similar to humans perception. A more detailed explanation about

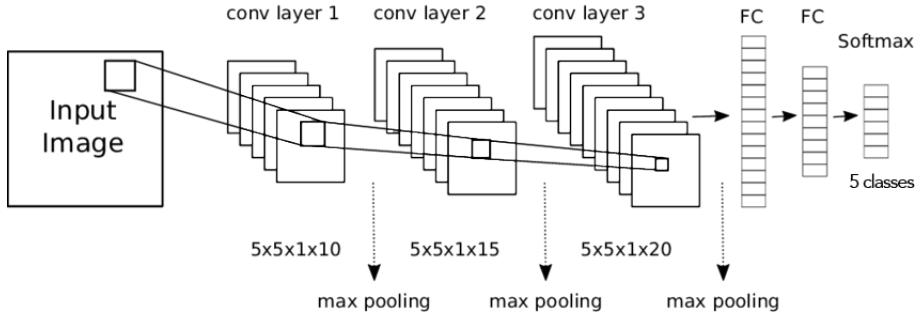


Fig. 3.5. Representation of a typical CNN model [80]

this can be found in appendix D. It has been used in plenty of works with CNN and also combined with other techniques [42] [81] [82].

3.2.2. Visual Geometry Group (VGG) model

CNN are usually able to achieve really good results and even improve human skills on Computer Vision tasks, for example, on recognizing object in an image. With the exponential growth of the researching works on this field, some challenges have appeared so as to promote the creation of new systems and test their efficiency and results. This is the case of the ImageNet Large Scale Visual Recognition Challenge (ILSVRC), based on the database of the same name, ImageNet. As a solution for the proposed exercise, the investigators from the Visual Geometry Group (VGG) in the University of Oxford implemented a new system achieving the first position and winning the challenge in 2014 [83]. The work they proposed consists in a study of the depth in a Convolutional Network architecture and how this can affect to the accuracy on the goal of large-scale image recognition [84]. To try this, it was necessary to increase the number of layers in the network, which was viable due to use a small size of convolutional filters in all of them.

For the training step of their system, they used an input image with standard size of 224×224 in RGB format. The principles to build the architecture are detailed below:

- The input image crosses a bunch of convolutional layers in which the kernel has a size of 3×3 .
- The convolutional stride has a value of 1 pixel.
- The padding is fixed to 1 pixel, so the dimensions of the input do not change during the convolution.
- Max-pooling is also included with a window size of 2×2 and a stride value of 2 pixels.
- Two Fully-Connected (FC) layers with 4096 channels after all the Convolutional Network layers.

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224×224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 3.2. VGG CONVNET CONFIGURATIONS.

- One FC layer with 1000 channels to perform the ILSVRC classification.
- A soft-max layer is employed for the final layer.
- All hidden Convolutional Network layers use the non-linear function ReLU as activation function.

All the designs that the creators came up with are based on these initial guidelines, except from just one case where Local Response Normalisation (LRN) is applied, previously mentioned in subsection 3.2.1. They just differ from each other on the number of layers, starting with 11 the first approach and ending with 19 the last one. The different architectures are specified in the table 3.2 [84] and are ordered from A to E.

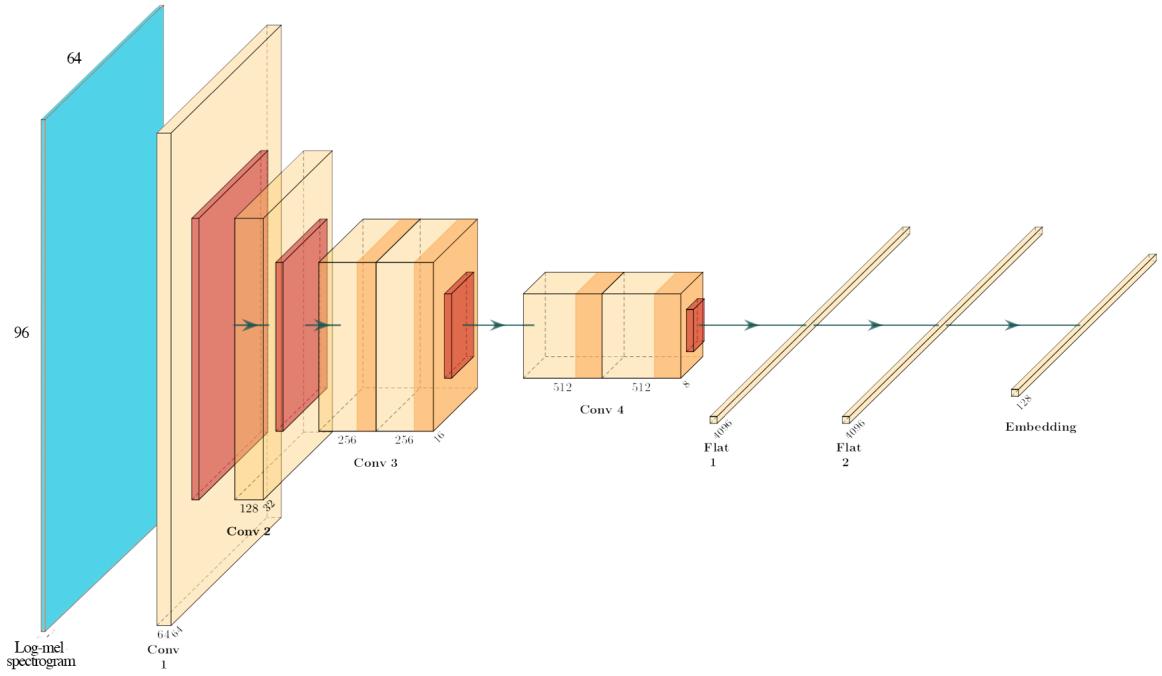


Fig. 3.6. VGGish architecture

3.2.3. VGGish model

The model we used in our project for feature extraction presents a configuration with principles similar to the ones explained in the previous subsection 3.2.2 but with slight changes that their developers have included in order to adapt it to the audio approach. As previously mentioned in subsection 3.1.4, the implementation of this model has been obtained from the website of the database Audio Set, in which a code repository is public available [85] and all the information about the implementation is included.

The architecture is based in the configuration A from table 3.2 with 11 weights. The differences respect to the original network are listed below:

- The input size was changed from 224×224 to 96×64 because of the log-mel spectrogram audio inputs.
- They built the implementation with just four groups of convolutional and max-pool layers, instead of five.
- For the last FC layer, they decided to build it just with 128 channels since it is the layer that compacts the final embedding and defines its dimensions.
- Also the final *Softmax* layer is not used.

Figure 3.6 shows how the configuration of the final VGGish model looks.

Input stage

Before passing the data through the whole CNN, a preprocessing stage has been included by the developers in which the input audio will suffer some transformations.

1. In the first place, after loading the input audio, the sample rate and the number of channels are checked to be 16 kHz and monochannel, otherwise the file is transformed to satisfy these conditions.
2. When the data is ready, the next step is computing the spectrogram. For this, it is necessary to calculate first the Short-Time Fourier Transform (STFT) of the signal. The operation is performed by using a sliding Hanning window of 25 ms with a hop size of 10 ms. As a result, just the magnitudes that correspond to the positive frequencies are retained, since the ones corresponding to the negative part of the spectrogram are redundant.
3. Transforming the spectrogram to Mel scale is what follows. To do so, they compute all mel frequency bins that are going to modify the values of the spectrogram in frequency domain by following a hertz to mel conversion formula, that can be found in appendix D. The result is a mel spectrogram from 125 Hz to 7,500 Hz divided in 64 bins.
4. Then, the log-mel spectrogram is calculated by doing the logarithm of the previous result plus a small offset value to avoid the logarithm of 0.
5. As a final step, they compute a framing operation over the log-mel spectrogram. The resulting are non-overlapping examples of 0.96 seconds, in which 64 mel bands and 96 frames are contained, each frame with a duration of 10 ms.

Therefore, the result obtained is an ensemble of 10 frames, approximately one per second, each of them with size 96×64 , i.e, 96 frames and 64 mel bands. It is worth to mention that this result size is obtained with the default parameters the system code is provided. These could be changed in order to obtain embeddings with a higher number of frames and an overlapping phase between them. Since, for our work, a comparison of the two formats was done, which can be found in subsection 4.2.3, and the already extracted embeddings in *.tfrecord* format follow the standard shape, we have not done more research in this aspect.

Embedding stage

Once the initial processing part is done and the log-mel spectrogram matrix is computed and divided into the desired number of frames, it is used as input data for the VGGish CNN. After all the computations inside the network, each example is converted into an embedding of size 128 giving a result of one of these per second of the original audio file.

Postprocessing stage

As final step, they performed some post-processing of the resulting embeddings. A Principal Components Analysis (PCA) [86] transformation is done joint with a whitening process. Also, a quantization to 8 bit for each embedding element. All these actions are computed with the purpose of making the final output compatible with the embeddings obtained from the YouTube-8M database.

3.3. Methodology

In this section, a theoretical explanation of the different algorithms and methods used in the experiments explained in chapter 5 is included, apart from the CNN, which has been already explained in subsection 3.2.1. More details about the implementation of the four techniques and how they play a roll in the system are included in subsection 4.3 and chapter 5.

3.3.1. Synthetic Minority Over-sampling Technique (SMOTE)

As previously mentioned in 2.2.1, SMOTE consists of an over-sampling method whose function is to generate synthetic samples in order to populate the minority class in a classification problem by making use of the proper observations that are already contained in the data [87].

Essentially, this technique works by choosing samples from the less populated class that are represented in the features space, placing a line between two of these samples and then selecting a point inside the line at a random position. First, an arbitrary sample is picked and then a number of k closest neighbours in the feature space are selected from the same class. One of the neighbours are randomly chosen and a synthetic sample is generated [88]. The number of neighbours k can vary depending on the amount of over-sampling needed. The standard value which is the one proposed in the publication paper is $k = 5$ [87].

For the generation of the synthetic samples, the distance between the feature vector first arbitrarily selected and one of the nearest neighbours is computed. Then, the resulting value is scaled by multiplying it by a random number that lies in the uniform distribution $U(0, 1)$. Finally, this new sample, S , is added to the initial feature vector and placed in between these two observations in the feature space [87]. S can be defined as follows:

$$S = x + u \cdot (x^R - x)$$

where x is the arbitrary selected sample the first time, u is the random value from $U(0, 1)$ and x^R is one of the closest neighbours [89]. In figure 3.7, a graphical representation of the algorithm is shown. In this case, X is the original sample for which the closest neighbours

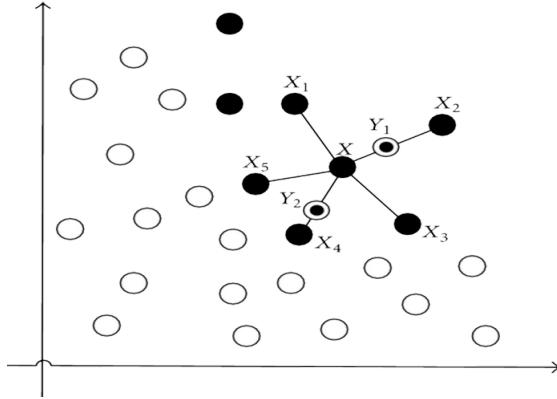


Fig. 3.7. SMOTE algorithm [90]

are represented as X_i where $i = 1, 2, 3, 4, 5$. The generated samples are Y_1 and Y_2 that are placed in the line segments.

The common way of applying this technique is in the train set before performing the fit of the model. Usually, an undersampling is performed in the most populated classes and, then, the data augmentation is done for the minority ones [88]. Our implementation follows this way of acting and it is explained in section 5.1.

3.3.2. Support Vector Machine (SVM)

Considering instances of data with their respective ground-truth labels, this classifier finds the frontier defined as a hyperplane that maximizes the distance between the two observations closest to each other that belong to different classes and are more likely to be misclassified. These two points are called support vectors [36]. In order to find this hyperplane, the data is mapped into a higher dimensional space so it can be more separable. This is done by applying a kernel function that eases the process of finding a proper mapping function. Some of the most used ones are the polynomial, Gaussian, linear, sigmoid and RBF [91].

This type of method is characterized for being a binary classifier. However, its implementation has been transported to multiclass problems by combining various SVMs in order to establish a decision criterion so the whole system can chose among Q categories. This can be done by following either one-versus-all or one-versus-one approach. The former addresses the problem training Q classifiers so as to differentiate between data from one class and data from the other $Q - 1$ classes. The latter treats the problem by training $\frac{Q(Q-1)}{2}$ SVMs to distinguish among all possible combinations of categories. In any of the approaches, the classification of a certain observation is performed by computing the distance between the sample data and the hyperplane that defines the frontier [9].

The one-versus-one is the one implemented by C-Support Vector Classification (SVC) from Scikit-learn library and is the one we will use in one of the experiments explained in section 5.2.1. In figure 3.8 it is shown a representation of the decision boundaries

that could be designed by other type of methods against the hyperplane a SVM classifier would draw.

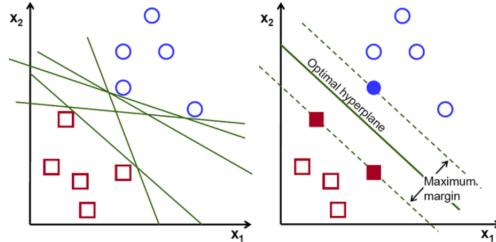


Fig. 3.8. Visualization of a hyperplane set by SVM and other decision boundaries [92]

3.3.3. Recurrent Neural Network and Long Short Term Memory

Human learning does not happen at each moment independently in the sense that every time something new is learned, it depends on the previous knowledge in order to interpret it. This can also be explained by saying that our understanding is persistence.

Traditionally, Artificial Neural Networks have not been able to act in this way. They cannot use time as a property to infer conclusions or predictions from the previous instances that belong to a sequence of data. To address this task, Recurrent Neural Network (RNN) have been developed, so they can take information and make it persistent. Figure 3.9 shows a block of a RNN in which the network, A , is fed with an input x_t that is actually a sequence of data. The output is the value h_t . An inner loop takes the information from the outcome and pass it to the input for the next learning step. This can be easily seen in the unrolled part of the right, how the information flows from one element to the next one [93].

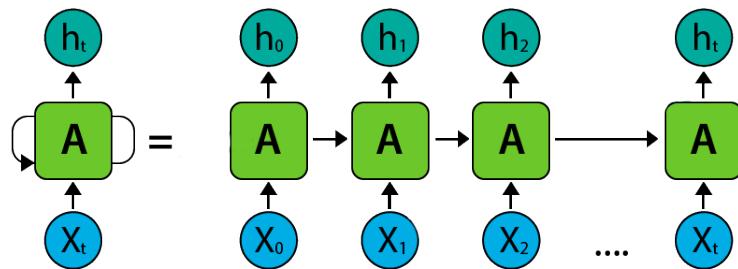


Fig. 3.9. Scheme of a recurrent neural network about how the loop works

This type of networks has became the standard when dealing with sequential data. They have been applied in many tasks, such as language modelling, speech recognition, translation, etc. However, the most simple approach present a problem when the information that must be considering when long-term dependencies. For example, if the task consists of predicting a new instance based on the previous ones in a not really long sequence in which the dependency resides on close instances, normal RNN can be used and work in the way explained above. When the distance inside the sequence between the

predicted element and the ones with the important information that this new one depends on is too large, the network cannot learn this connection [93].

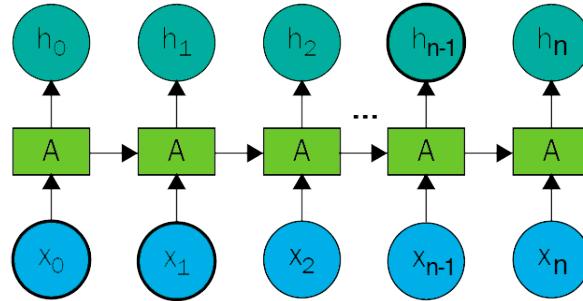


Fig. 3.10. Example of long-term dependencies situation. If the predicted output h_{n-1} depends on x_0 and x_1 , the vanishing gradient problem may appear.

This happens due to the value of the gradients becomes very small during the back-propagation in the learning process. Once the loss function for the new predicted output is calculated, this must be propagated through the rest of the network in order to update the weights. In a typical neural network model, the weights that get updated are those from the hidden layer directly previous to the output one, but, since the data is a sequence, the ones from all the earlier layers must be updated as well. The actual problem appears when renovating the value of the weights through time, i.e. updating the weights used to connect the hidden layers in the unrolled temporal loop. So, if the value of the gradient becomes very small the updating process start to be null when finding the new values of the weights and the model stop learning. This is known as the vanishing gradient problem [94].

As a solution, a new type of RNN was developed and named as Long Short Term Memory (LSTM) networks. These can be defined as a recurrent network that can learn long-term dependencies, so the vanishing gradient problem is not an issue for these models. Nowadays, they are widely used in several fields, since they can remember information along long periods [93]. The main structure along the model is the same shown above in figure 3.9. The difference from the way of working in a common RNN is the inner architecture in the neural network. In LSTM, the A modules are defined as shown in figure 3.11. The notation used is as follows: the rectangles denotes a neural network layer, the circle shapes refer to point-wise operations and the arrows means vector transferring.

The main concept of a LSTM network resides on its cell state and the different gates. The cell state, represented in the diagram by the upper arrow, is the one in charge of transferring information all the way through the chain of modules. It can be thought as the memory of the whole network. Its objective is to just carry information considered relevant for the model. This way is how the LSTM brings the information from the first layers to the last ones without suffering the vanishing gradient problem. The content of the cell state is modified by the point-wise operations that acts by following the gates criteria. These are in fact neural networks that has the goal of deciding which information

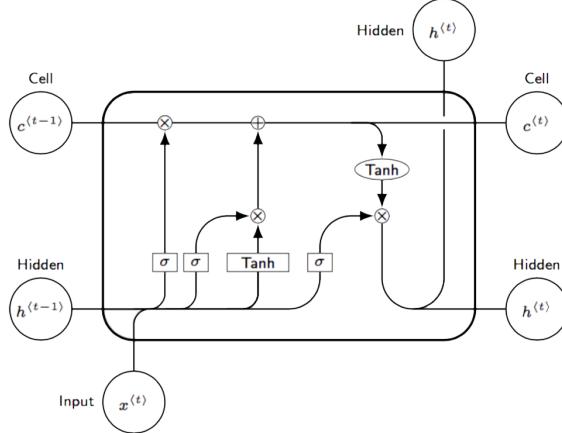


Fig. 3.11. LSTM module with the representation of the different operations that take place inside of it

can be treated as relevant and so added to cell state. This process can be understood as a learning/forgetting stage [95].

In order to understanding how this process works, we are going to address the interaction of the different gates in the learning stage.

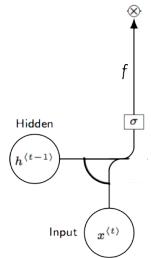


Fig. 3.12. Forget gate

The first step is performed by what is called the *forget gate* layer. It must decide which information must be included in the cell state that travels through the time chain. This is why the layer contains a sigmoid function that outputs values between 0 and 1, in which 0 means dropping the information in the previous cell state and 1 implies to keep it. As shown in figure 3.12, this is done by paying attention to the hidden state which is the output of the previous module, $h^{<t-1>}$, and the input of the current one, $x^{<t>}$ [95]. The resulting function f is defined as follows, where W_f are the weights of the RNN and b_f the bias.

$$f = \sigma(W_f[h^{<t-1>}, x^{<t>}] + b_f)$$

Next, the decision about what information from the current input $x^{<t>}$ must be taken to the cell state is performed. The gate that computes this operation receives the name of *input gate*. This part can be divided into two different steps. First, the previous hidden state, $x^{<t-1>}$, and the current input feed a RNN with a sigmoid for the activation function, which takes the decision of which values must be updated by converting them into a range limited by 0 and 1. This result denotes the importance of the value, being 0 non-important and 1, important.

Then, also $h^{<t-1>}$ and $x^{<t>}$ are used as input for the tanh layer so the values are mapped between -1 and 1. This last activation function is included in order to avoid the vanishing

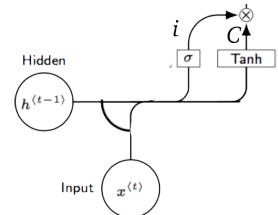


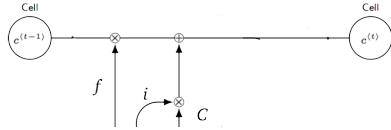
Fig. 3.13. Input gate

gradient problem previously mentioned, since a function whose second derivative takes more time to tend to zero is needed. Finally, both outputs are multiplied before the updating process in the cell state [95]. The equations for the output of each RNN module are included below, where W and b denotes the weights and bias for each layer [93].

$$i = \sigma(W_i[h^{<t-1>}, x^{<t>}] + b_i)$$

$$C = \tanh(W_C[h^{<t-1>}, x^{<t>}] + b_C)$$

In the step 3, the output of the multiplication from the input gate layer and the output f of the sigmoid function from the forget gate, modify the old cell state from the previous module, $c^{<t-1>}$, so it can be updated. For the outcome of the step 1, a point-wise multiplication is performed, $f * C^{<t-1>}$. Then, to the output of this operation is added the resulting product of the input gate, $i * C$. This two operations can be translated as a forgetting and a learning stage. First, the values the forget gate decided to remove from the cell state are actually forgotten. Then, it must learn the new values belonging to the current input also weighted by their importance, what denotes how much the state values are going to be updated [93]. The new cell state expression c^t is included below and the process is shown in figure 3.14.



$$c^{<t>} = f * C^{<t-1>} + i * C$$

Finally, the last step of the module corresponds to the *output gate*. This has the function of deciding how the next hidden state is going to be like. First, $x^{<t>}$ and $h^{<t-1>}$ are used as input for a layer with a sigmoid activation function and the output o is obtained. Also, the cell state $c^{<t>}$ feeds a point-wise tanh so its values are clipped between -1 and 1. Both outputs are multiplied. The o result decides what values form the current input must be kept in the future hidden state, $h^{<t>}$. The actual output of the module is the hidden state. It also transports the relevant information updated together with the cell state to the next network. The outputs o and h^t can be defined as shown below. Also, in figure 3.16, the stage of this last part is included.

$$o = \sigma(W_o[h^{<t-1>}, x^{<t>}] + b_o)$$

$$h^{<t>} = o * \tanh(C)$$

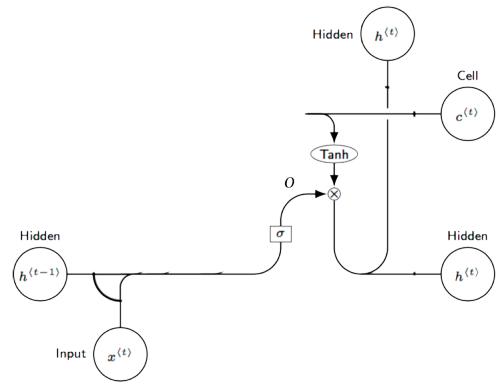


Fig. 3.16. Output gate

4. OUR APPROACH FOR ACOUSTIC VIOLENT DETECTION

For our final approach when implementing this system of acoustic gender-based violence detection, we first need to select the classes we are going to use for the classification so as to obtain the corresponding embedding features with VGGish, either from the audio data or from the *.tfrecord* files obtained from the database explained in section 3.1. Then, a supervised multiclass classification is performed for the different events and, finally, the soft output from this step is used as an input for a binary classifier so we can differentiate the audio from violent and non-violent.

In this chapter we want to make clear why we think that the VGGish network is a good option for our problem in 4.1. Then, the procedure when starting to prepare the data and to extract the feature embeddings with this model is included in 4.2. Also, the demonstration of the similarity between both types of embeddings previously mentioned in 3.1.4 is described in 4.2.3. Finally, the last section 4.3 includes an explanation of the whole system that we use for the experiments explained in chapter 5.

4.1. Why VGGish

For our goal, one of the toughest tasks consisted on the selection of data that properly adapted to our problem and its preparation so as to obtain features that allowed us to characterize every acoustic event from a violent point of view. Since our first efforts of finding an available dataset containing enough violent scenes was driving us to a dead end, we decided to take advantage of a huge database which let us rethink the standpoint about how we were going to address the problem. As it was mentioned in subsection 2.3.2, one of the main questions was how to define the term violence for each victim depending on her certain situation. After finding the Audio Set database, previously explained in section 3.1, with all its variety of samples, we had a wide range of audio data to work with. This is how we came up with the system explained below in 4.2.1.

At this point, not only we had an idea but we had already found a data resource to start with. However, the issue was related to what kind of features could be extracted in order to categorize events from different nature with a unique violent label. An acceptable conception of the term violence could be expected to cover all kinds of impulsive events such as hits, smashes, gunshots, yells, etc. Also, we would like to introduce sounds that were likely to happen in a domestic environment within a tense atmosphere as children crying, dog barking or glass breaking. We also wanted to take into account the possibility of including other cases not usually considered violent a priori. For example, the sound of keys jangling or the noise produces by a shaver machine. These situations may be too particular and just would be present in few uses, but this is how we understand the problem given the advice of the social sciences experts in the UC3M4Safety team.

So, our first intention was to apply some audio processing techniques to extract low-level features, as the ones previously explained in 2.2. Even though there are plenty of previous works and a lot of tools to work in this way, it was not sure which path should we take in order to decide what features better fitted our task. Apart from this, since the database had such a big volume it would have supposed an enormous cost of time to compute features every time we wanted to try new type of categories. Moving on, by following the advances on finding new level features already mentioned in 2.2, we decided to investigate new methods of extraction based on the use of Neural Network models. Nevertheless, even thought the features obtained in this case had been more appropriated, the time consumption of training a big model was one of the aspects that did not totally convince us.

The previous selection of Audio Set as our dataset allowed us to get to know the VGGish model proposed by Google researchers for feature extraction. This system loads the parameters already learnt from training with another huge dataset as YouTube-8M. This is possible due to apply transfer learning idea, explained above in subsection 2.2.2, that consists of leveraging features or weights extracted from certain models and use them in simpler ways for different tasks [48], so all the computational cost and training time is not a problem any longer. Finally, we decided to put in practice this pre trained embedding extractor by loading the given parameters so as to obtain our final input representations.

4.2. Our approach

In order to start describing our approach, we will first explain in 4.2.1 how we obtained the input data to work with by using the resources previously explained in 3.1 and 3.2. Then, we comment some considerations when extracting the feature embeddings in 4.2.2 so as to finally explain in 4.2.3 the difference between the two types of data.

4.2.1. Input data

Different phases took place when trying to obtain all the necessary data from the YouTube videos specified in the database. We will explain them from the first step of deciding which classes better fit our problem to the last part in which the desired embeddings are achieved.

Violent classes

In subsection 2.3.2 we mentioned our idea about giving the victim a choice of defining her own perception of violence, so the final implementation can adapt to her situation in a better way. To do so, we have taken advantage of the ontology provided by the Audio Set creators we explained in subsection 3.1.2.

Our system has been implemented based on the idea of using the parent-children relationship among the different nodes to ease the process of selecting among the 632 classes. It must go through all the branches so as to offer the victim the possibility of choosing any of the audio event categories. However, instead of consider each label individually, this starts the way from the parent classes down to the children ones.

Let's say we begin from the class "Human sounds" that is the top level of all sounds emitted by humans contained in the dataset. The system will ask the user if she wants to advance in that direction, i.e., to go across the branches that belong to that part of the tree. If the answer is positive, it would go for the next class, that in this case it would be "Human Speech". It will advance this way until there were no children nodes in the actual class. When this happens, the user will be asked to add the label to the record of *violent classes*, that will be saved in a text file so they can be read by other parts of the model further on. If the user does not want to go deeper, she will be asked to add the current class to the record. If the answer is "No", then they system will jump to the next sibling category. The corresponding flowchart is shown in figure 4.1.

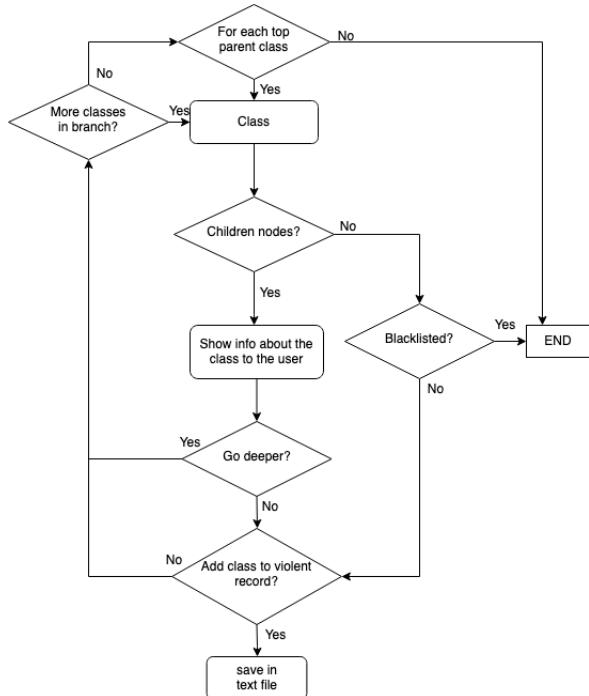


Fig. 4.1. Flowchart about selecting violent classes
The corresponding flowchart is shown in figure 4.1.

This way of flowing through the different classes allows to skip those parts that are not related to our problem. For example, as we can see in figure 3.1, one of the top parent labels is "Natural sounds", that relates to sounds from weather phenomena. In most of the cases, these classes will not be selected so the whole branch can be skipped.

Downloading videos

The following step consists in achieving audio files that belong to the chosen labels. For this purpose, we have made use of the .csv files that were explained in subsection 3.1.4. For each included video, we took its ID and build the corresponding YouTube URL. Once downloaded, we trimmed the file considering the onset and offset and, finally, converted it to audio format (.wav). In our script, we can pass as a parameter the identifier of the desired classes in comma-separated format and either the number of videos per class for a balanced set or a total number of downloads for an unbalanced set. However, there might be some errors when obtaining all the data: the two most common cases are due to lack

of enough videos of the desired type in the dataset or because the video is not available anymore on YouTube. When this happens, a message is shown to the user.

It is also worth mentioning that throughout the developing of this downloading task, a script has been coded to achieve the whole dataset in both formats, video and audio, for future works. We are not going to specify anything else since this feature was not finally used.

4.2.2. Extracting embeddings

At this point, all the desired data has been already downloaded to extract the embedded features that will be used to train the model. For this part, we have used the VGGish network explained in subsection 3.2.3. Since the audio files duration is usually 10 s, and the embedding extractor gives as a result a vector of size 1×128 for each second, we will obtain a 10×128 feature matrix composed by values within the range 0 - 255. Therefore, our input feature matrix will have a size of $(\text{number of audios}) \times 10 \times 128$. The corresponding labels will be stored in a vector of size $(\text{number of samples})$.

There are some points about the data obtained in this step that should be commented. One of them is about the length of the audio files. As it was previously indicated in 3.1.1, most of them last 10 s because the creators decided to set this duration for the audio events, but this can change if the video is originally shorter. For these situations, since the model is configured to have an input of 10×128 , the embedding matrix of the shorter clip must be padded with zero-rows to achieve the required dimensions. Even though this does not happen very often, there might be some silent segments that will be labelled with the category of the rest of the audio. A solution for this problem explained later in 5.1 was implemented.

The other case is related to what was explained in 3.1.4. For the recent explanation of how to extract the input features from the audio files we have utilized the first manner of accessing the data, i.e., by reading the .csv files with the videos information before downloading. There is a second option of using the already extracted embedding features. However, these do not look exactly the same when comparing them to the ones obtained by running the VGGish with the downloaded videos. This difference is because the implementation of the given code differs from their internal production system in computing issues such as underlying libraries in the installation of VGGish and hardware equipment. In spite of this, the result in classification tasks are expected to be equivalent. In order to prove this assumption, we decided to try a small system with both kinds of data. The reason for this check is that although, for this work, the development of the methodologies could be made by using the first and faster manner of accessing the data, for the real implementation in the pendant device within the EMPATIA project (out of our scope) only the second way will be feasible.

4.2.3. Assessing the differences between the two types of data access

In order to check what is explained above in 4.2.2, we have decided to run a little experiment in which a toy classification is performed. Also, we wanted to visualize the different features to check if we could appreciate common patterns by using the t-distributed Stochastic Neighbor Embedding (t-SNE) algorithm, which will be explained later in this section.

Our first step consisted in determining a subset extracted from the original dataset. We thought about choosing for this small application a set composed by three classes that could be considered violent and other three that were non-violent. Apart from this, we paid attention to the number of samples per category to pick some classes over others. Finally, we ended up picking up the labels detailed in table 4.1 and a number of 80 samples for each of them, which led us to a total of 480.

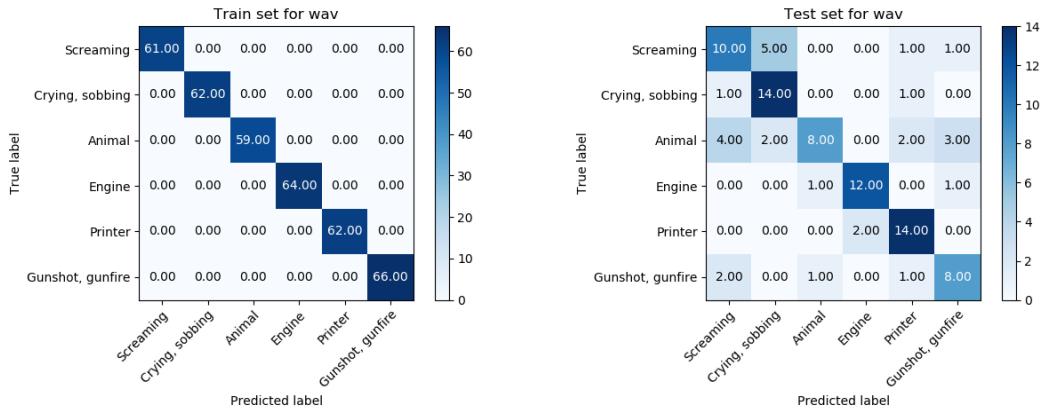
Class	Description
Screaming	A sharp, high-pitched human vocalization; often an instinctive action indicating fear, pain, surprise, joy, anger, etc. Verbal content is absent or overwhelmed, unlike Shout and Yell.
Crying, sobbing	Sound associated with the shedding of tears in response to an emotional state, arising from slow but erratic inhalation, occasional instances of breath holding and muscular tremor.
Gunshot, gunfire	The sound of the discharge of a firearm, or multiple such discharges.
Animal	All sound produced by the bodies and actions of non-human animals.
Engine	The sound of a machine designed to produce mechanical energy. Combustion engines burn a fuel to create heat, which then creates a force. Electric motors convert electrical energy into mechanical motion. Other classes of engines include pneumatic motors and clockwork motors.
Printer	Sounds of a computer peripheral which makes a persistent human readable representation of graphics or text on paper or similar physical media.

Table 4.1. CHOSEN CLASSES FOR A SMALL CLASSIFICATION.
*SCREAMING, CRYING, SOBING AND GUNSHOT, GUNFIRE ARE
 CONSIDERED AS THE VIOLENT ONES.*

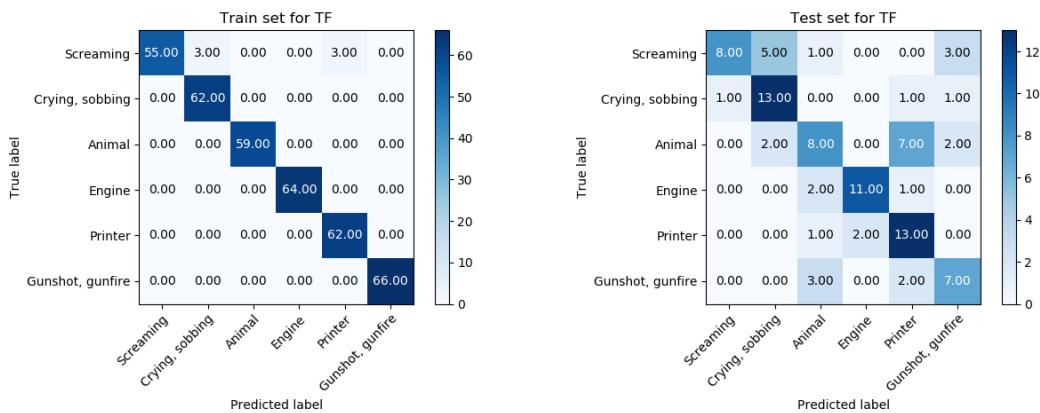
For the classification task, we decided to create a small CNN composed by few layers. Since our input data are matrices of shape 10×128 , these were treated as images so the

model was built with layers that perform spatial convolution [96]. The architecture is detailed in figure 4.3. We used a small kernel size of 3×3 , zero padding so as not to change the shape of the output and a activation of function ReLU. Two dense⁵ layers are added at the end, first one with also ReLU as activation function and the second one with softmax to perform the classification, and as many filters as the number of classes.

In order to measure the results, since our subset is balanced, we could have evaluated our model by computing the accuracy and the confusion matrix. More information about these metrics can be found in the appendix A. In figure 4.2, the four confusion matrices corresponding to the training and test phases for both types of data are shown. Also, in table 4.2, it is included the accuracy for each case. The results are more accurate when training with the embeddings extracted directly from the audio files we downloaded. When obtaining them from the TensorFlow files, they value of the metrics indicate a worse performance. However, we opted for this second manner because the performance difference is not too big and it requires much less computational cost and time loss.



(a) Confusion matrices when embeddings are extracted from audio files



(b) Confusion matrices when embeddings are taken from .tfrecord files

Fig. 4.2. Confusion matrices

⁵Dense is a also a common form to refer to Fully-Connected layers. Since each neuron is fed with an input from all the neurons in the previous layer we can say that they are densely connected [97]

format / subset	Train	Test
Audio file	1.00	0.70
.tfrecord	0.98	0.64

Table 4.2. ACCURACY VALUES FOR AUDIO AND *.TFRECORD* FILES

For the training phases, it can be appreciated that there may be an overfitting since the accuracy is perfect. This means that the NN stop improving its capacity of learning how to solve the problem in a certain moment of the training task. Instead, it does learn some behaviour pattern that the training data follows. This impacts negatively in the model since the new data that the system will have to learn from will look different and will not follow these same rules [43]. In spite of this result, we did not give it so much importance since we just wanted to do a quick check of the similarity between the two types of data which can be appreciated due to the similarity of both metrics results.

Dimensionality reduction for visualization

Apart from the classification exercise, we wanted to see if by plotting the data samples we were able to identify or appreciate some common patterns. In our problem, each of our samples is characterized by a matrix of features with 128 columns, which means that we were working with data belonging to a high-dimensional space. Visualizing this type of data has always been a case of study for many different fields. Plenty of methods have been published so as to find a solution for this task. Some of the most accepted consist on reducing the dimensionality of the data so this can be transformed from the high-dimensional space to a lower one and can be visualize in a common scatter plot of 2D or 3D. These techniques rely on the idea that within a multivariate sample denoted as $x_i = [x_{i1}, \dots, x_{in}]^T$ and considered to be a point that corresponds to a $n - \text{dimensionality}$ space, a $d - \text{dimensionality}$ space can be found, so that $d < n$. If this is possible, then the observations can be transformed to this lower dimensional space d without any loses [98].

One of the most common and well-known dimensionality reduction methods is the one known as Principal Components Analysis (PCA). This follows the idea previously explained. It specifically aims at extracting the *important* information from the original data and transform them into a set formed by orthogonal variables which are actually

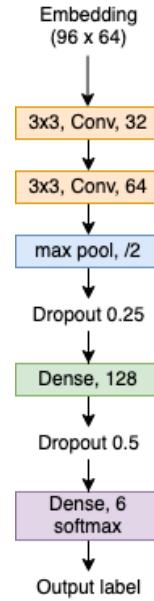


Fig. 4.3. Architecture to see how the different embeddings work

known as principal components. This is done by multiplying the matrix data X by a projection matrix Q that contains the coefficients of the linear combinations that allow the conversion. The projections must be orthogonal to each other and they represent the data maximum variance in descending order, being the first component the one with largest variance [86]. In fact, each of the projections correspond to an eigenvector in descending order following the value of the eigenvalue. So, the first component will be the eigenvector with the highest eigenvalue. It has been proved to be one of the most appropriated options in dimensionality reduction nowadays. Its use is completely accepted and it is implemented in many famous software libraries. However, it presents some limitations. One of them consists on just considering linear combinations of the original data. When the relation is non-linear, a dimensionality reduction with this technique may result in a loss of information [99].

For these cases, other methods have been developed such as t-distributed Stochastic Neighbor Embedding (t-SNE). This algorithm appears as an extension of the previously published Stochastic Neighbor Embedding (SNE) [100]. Both are based on the same idea of a new way of measuring the similarity between samples. Instead of comparing two observations, let's call them x_i and y_j , by computing the euclidean between them, this is done by calculating the conditional probability $p_{j|i}$ of x_j being picked as a neighbour of x_i considering that the samples belong to a Gaussian distribution centered at x_i . This depends on how far the samples are from each other, i.e., it is high when they are close and minimum when there are totally separated [101]. Apart from this, two analogous observations are created in the subspace of lower-dimensionality, y_i and y_j , and conditional probability $q_{j|i}$ is computed in this situation. It is important that, for y_i and y_j to be faithful representations of x_i and x_j , both conditional probabilities must be equal.

In order to calculate the probabilities, a crucial factor is the variance of the Gaussian distribution. There is not a unique valid choice for this parameter, so t-SNE performs a binary search so as to find the optimal one [99]. This is also influenced by what is called the perplexity. This can be defined as an assumption of the number of adjacent neighbours for each point. It is a value that is fixed by the user, but it usually belongs to a range from 5 to 50 [101].

There are several considerations that should be known before looking at a representation of data from this algorithm [102]. Actually, it is not an easy task to understand this kind of plots, since the distance between points in the new subspace are not related to the real euclidean distance, which has been denoted as "The Crowding Problem" [101]. This means that the groups cannot be interpreted as real collections of data in the original dimension. However, in order not to misunderstand the data distribution, a couple of visualizations varying the parameters usually tends to be done so that the conclusions can be based on more than one result. In figures 4.4 and 4.5, we show ten t-SNE outputs, five for each type of data for diverse values of perplexity from 2 to 50, respectively.

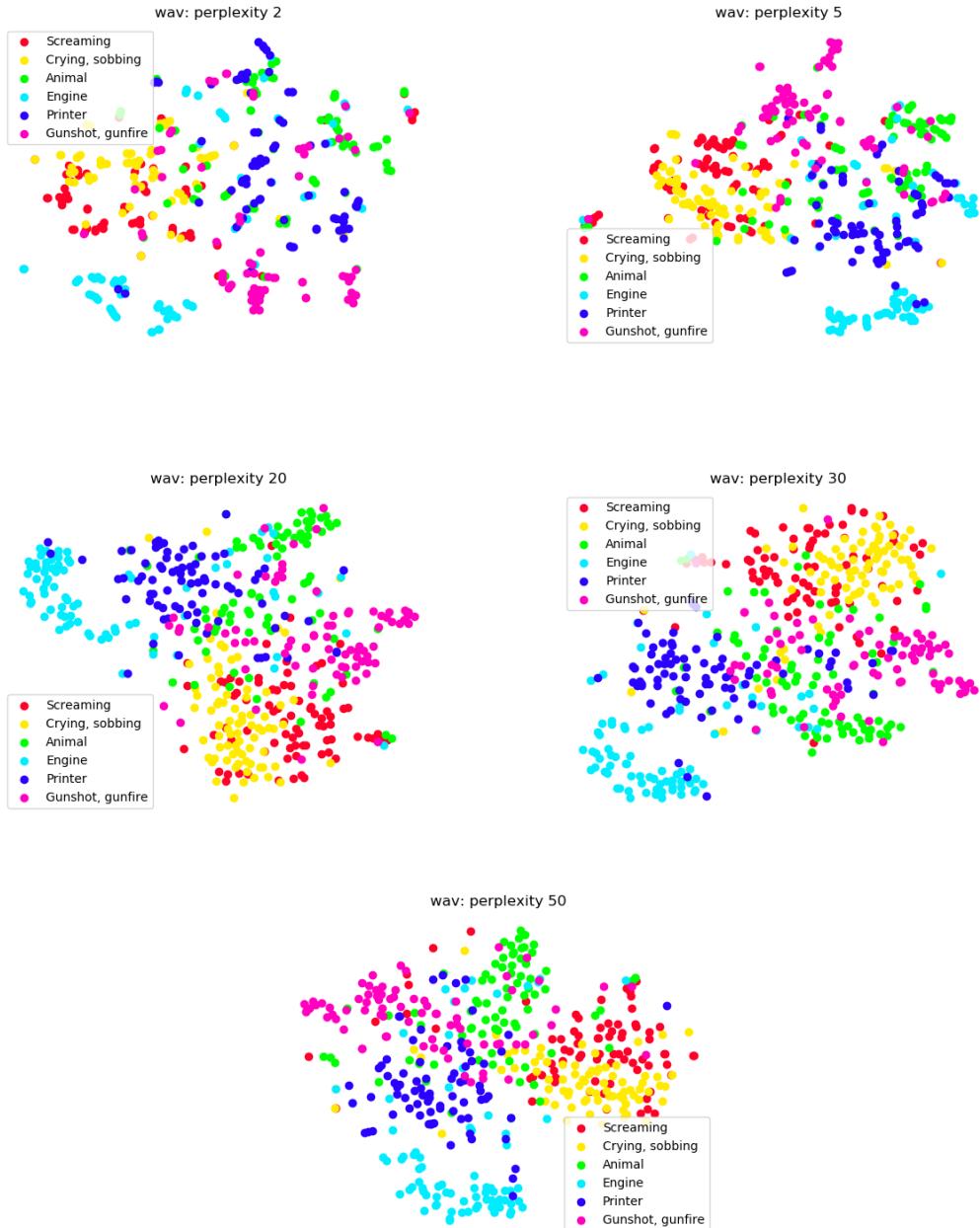


Fig. 4.4. t-SNE results for wav format with a legend that shows the labels of the corresponding samples in the original 128D space

The selection of the perplexity value depends on the number of observations per class [102]. Since our subset has 80 samples for each category, we can consider that a proper value is 20 or 30. As we can see, a boundary cannot be extracted among the different labels, but we can see some grouping patterns in the data that correspond to the original labelling and helps us to confirm the similarity of the two different given types. It is true that this method should never be used as an algorithm for clustering itself, but it is a good resource as a backing strategy to other results, as in this case.



Fig. 4.5. t-SNE results for *.tfrecord* format with a legend that shows the labels of the corresponding samples in the original 128D space

4.3. Our system

In this section we are going to make a small explanation of the systems we have implemented. As explained at the beginning of this chapter, our main objective in this work is to perform a binary classification in which we can distinguish between violent and non-violent events. Before achieving this final result, a couple of steps take place from which we can make some pre-processing of the data and also extract some conclusions about the behaviour of the embeddings selected. More information for the implementation will be

found along chapter 5.

First, we need to select the classes we are going to work with. For that purpose we use the system implemented explained in subsection 4.2.1. After this, we have ended up with a couple of violent and non-violent classes that will be used in the classification task.

Then, due to the unbalanced nature of the database, some of these classes are more populated than others. We decided to solve this problem by using a data augmentation technique called SMOTE which generates samples for the minority classes in a synthetic way. A more detailed explanation can be found in 3.3.1. We actually did not have a very wide range of possibilities when choosing if a synthetic manner of creating new data was the best idea. Since we are not working with original audio files and we are just taking the already extracted embeddings provided in `.tfrecord` format, applying some of the techniques about deforming or transforming the original signal, as previously explained in 2.2.1, was not a possible option. SMOTE is a common technique for this task with a robust algorithm that seems a good election.

Then, after dividing the data in the corresponding sets for the classification, we wanted to see how was the behaviour of the embeddings in different ways. Then, the final decision was to perform a multiclass classification, as the small experiment described in 4.2.3, with different techniques and see how they work with this type of data. In the following list, the algorithms already explained theoretically in section 3.3 and 3.2.1 are included with the reason of why choosing them:

- We decided to use SVM since we wanted to use a classifier algorithm that differs from neural models. This technique has been used largely in the literature [36] and seems one of the most appropriated options when performing classification tasks. We also used this algorithm to perform the binary classification.
- Another option was to use the simple CNN model explained in figure 4.3. Since the data we are working with is a 2D matrix for every segment of audio, we wanted to check the result of a convolution operation around each segment. The number of channels indicated in this case is 1, so every instance is considered as a single input "image".
- As a last option, we wanted to try the concept of a LSTM, since this method takes as an advantage the fact of the input being a sequence. We could have used simpler RNN but we decided to try this method since the concept they present are mainly the same.

These three methods are used to perform a multi classification task for all the samples collected within the classes chosen. The output we get are the probabilities for each sample to belong to each class, i.e. the soft outputs. To evaluate the performance of this task, we convert them to hard values by choosing the higher probability. This will be also mentioned later in section 5.1. Actually, this is due to the softmax layer we include at

the end of CNN and LSTM. For the SVM, the outputs are naturally hard, we just need to indicate that we want to get soft ones in the prediction stage when declaring the model. So, depending on which classification task we are performing, we indicate hard or soft.

For the binary classification, we use the soft probabilities mentioned above as input features for the binary SVM. In this case, we get the hard outputs *0* or *1*, depending if the observation is non-violent or violent. In figure 4.6, a block diagram of the whole system is included. This shows the whole treatment of the data since being an audio recording until the output of the classification.

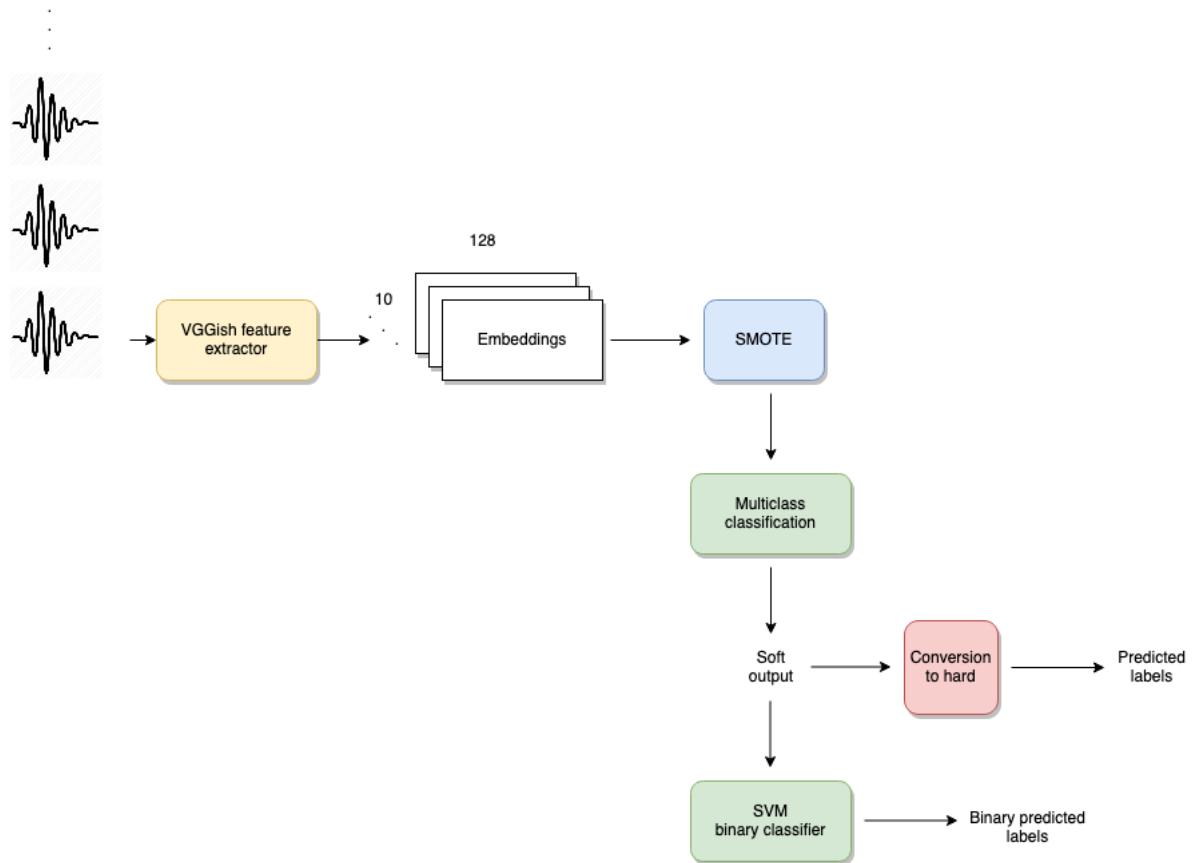


Fig. 4.6. Block diagram of the whole model including both types of classification. We do not extract the embeddings, but this is a diagram of the whole model considering that someone did extract them.

5. EXPERIMENTS

In this chapter we are going to explain the final models we chose for our work and the input data preparation to feed these models with. For the different solutions, we have used the algorithms previously explained in the methodology section 3.3. Below, it is presented a list including the 6 implementations compared:

- SVM classifier for multiclass classification
- SVM multiclass + SVM for a final binary classification
- LSTM for multiclass classification
- LSTM multiclass + SVM for a final binary classification
- CNN for multiclass classification
- CNN multiclass + SVM for a final binary classification

5.1. Input data preparation

For the proposed experiments, we have decided to build a small dataset with the embeddings extracted from the *.tfrecord* files that belong to 14 classes: half of them *violent* and the other half, *non-violent*. To find these, we first did a run on the simple user-interaction program that is explained in 4.2.1 as if we were gender-based violence victims. As we said, a total of 28 classes were selected. From this set, we picked 7 as the violent classes. The other 7 were chosen just by looking at the ontology to provide negative classes.

It is important to mention that this collection of data is composed by samples with just one label assigned. As explained before in 3.1, the Audio Set database is highly unbalanced. Some of the most appropriated classes to be considered violence are scarcely populated. When doing the selection of categories, apart from paying attention to the own meaning of the class, we also checked the number of samples. For the non-violent type, this was not a problem, since we took some of the most populated labels. However, in the violent case, we had to deal with the requirement of representing violence and also having enough observations. So, due to these limitations, we could not obtain a naturally balanced set. In table 5.1, the 14 selected labels are shown. In figure 5.1, a bar plot shows the number of samples obtained per class right after the selection.

Some initial preprocessing steps were performed before passing the data to the different models. As mentioned in 4.2.2, we had to deal with the zero-filling problem, which means that some of the rows of the embeddings matrices are completely zero because the

Violent	Non-violent
Baby cry, infant cry	Printer
Slap, smack	Music
Screaming	Speech
Machine gun	Vehicle
Breaking	Animal
Slam	Dishes, pots, and pans
Yell	Wind

Table 5.1. RELATION OF *VIOLENT* AND *NON-VIOLENT*
SELECTED CLASSES

duration of the original video is less than 10s. As a solution, we decided to substitute the zero numbers in all the data with the machine epsilon⁶ value.

However, in the two models that use SVM for the multiclass classification a different solution was proposed. This algorithm needs the input data matrix to be in the form [*number of samples* × *number of features*], which differs from the originally shape of our data, [*number of samples* × *number of seconds* × *number of features*]. For this reason, we decided to reshape our data to the required form which resulted in a matrix of shape [(*number of samples* × *number of seconds*) × *number of features*]. With this conversion, instead of working with full audio instances, the data samples became the seconds of those instances. In order to remove zero data, we first checked the amount of zero-rows in every class. Since it was not a very significant portion of the data, we took them out of the dataset.

Once we had our data with all non-zero values, we needed to convert the unbalanced set to balanced. First, for the classification task, we divided our data into train, validation and test subsets, using a 20% for last one. Then, we decided to exclude half of the samples from the most populated classes from the train and validation sets⁷. In figure 5.2, subfigure (a) shows the distribution of data per class in the dataset for the models that use SVM multiclass classifier. In subfigure (b), the distribution for the other four cases is shown. Then, we wanted to generate new data for those with less observations by applying the data augmentation technique SMOTE, explained in subsection 3.3.1, in order to generate samples until equalize the most populated category. Just the train and validation sets were subjected to this conversion process, leaving the test set with the original number of embeddings in their natural a priori state. The final shape of the input data for each of the experiments is shown in table 5.2.

⁶The machine epsilon value is considered the smallest value that satisfies $1 + \epsilon_{\text{match}} > 1$. It is the difference between one and the next closest number that is representable as a machine value [103]

⁷*Speech, Music, Vehicle* and *Animal*

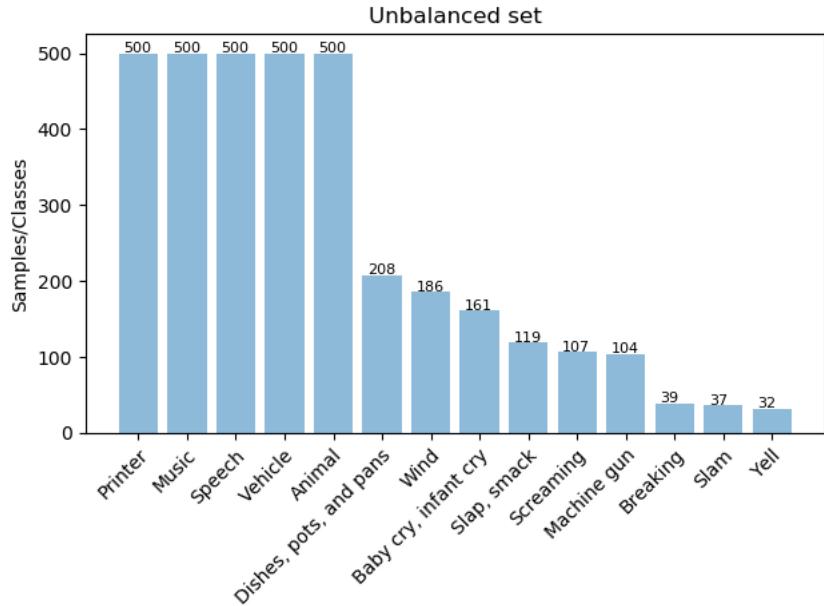
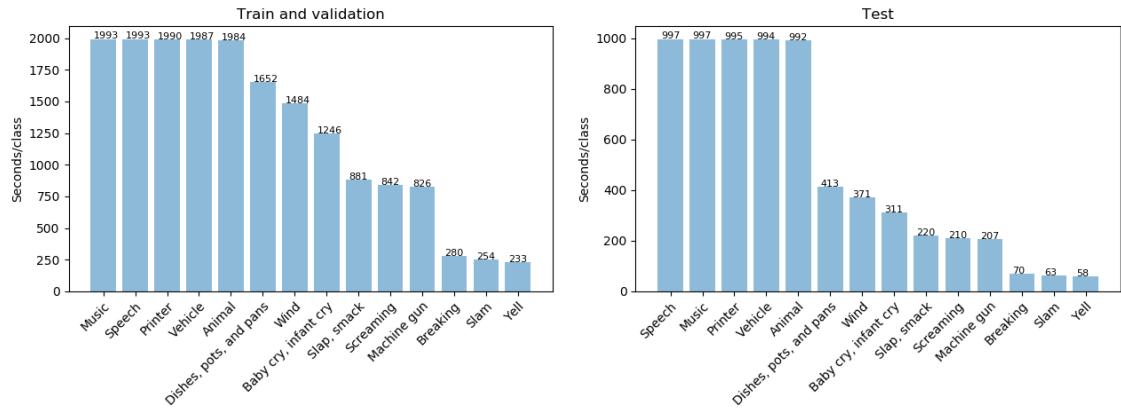
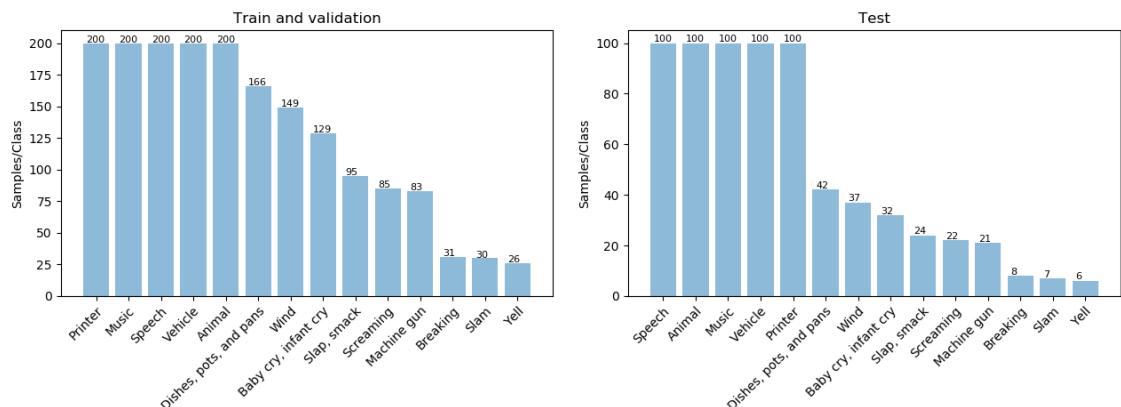


Fig. 5.1. Bar plot that shows the number of samples for each of the selected classes. Clearly, the violent categories are much less populated than the others, that did not have to much any semantic criteria



(a) This are the resulting subsets after downsampling the most populated classes and removing the zero-rows for the experiments that involve an SVM classifier



(b) This are the resulting sets after downsampling the train and validation sets. The test set remains the same after the split.

Fig. 5.2. Number of observations for train, validation and test subsets used in the different experiments 52

	SVM multi and SVM + SVM binary	LSTM multi and LSTM + SVM binary	CNN multi and CNN + SVM binary
Train and validation	17645 (1993 per class)	2800 (200 per class)	2800 (200 per class)
Test	6898	699	699

Table 5.2. TRAIN, VALIDATION AND TEST SET FOR DIFFERENT EXPERIMENTS

5.2. Implementations and results

For all the experiments, as mentioned above, the dataset was split by randomly selecting a 20% of the data for the testing procedure. The other part was divided into train and validation subsets with a k-Fold Cross Validation technique, with 10 folds, i.e. 10-fold cross-validation. A more detailed explanation about this resampling technique can be found in appendix B. The average and standard deviation of the results from the different folds were obtained for train and validation. Then, a final measurement was performed for the test set.

The way of checking the model performance was by finding the accuracy and confusion matrix, whose explanation can be found in the Appendix A. For the cross-validation procedure, a matrix with the average values⁸ was obtained and also one for the standard deviation,⁹ so the estimation of the model is represented for the corresponding subsets. Finally, a last evaluation of the model is performed by checking the accuracy just once with the test set. The testing step is performed by using the model for which the highest value of accuracy was obtained during the validation process.

Also, for the test set, in order to show an orientation of a binary solution performed with the original embeddings, the multiclass matrix result is adapted to a binary approach by building a new matrix that relates the number of predicted labels for the violent and non-violent classes. This way, we want to give an idea of how a binary classification directly performed with data labeled with just these two tags would be ended.

In order to compare the results for the three sets, an error bar plot is also shown for each implementation. This includes the average value of the accuracy for the train and validation set and the accuracy obtained the only time the model was evaluated with the test set. These are shown in a two decimals format for a better visualization. About the standard deviations, they are also included in the plot if they are high enough. Sometimes,

⁸In the average matrices, the results for each cell are shown just with one decimal in order to make a better and more comfortable visualization. However, the colour bar on the right of the plots must be taking into account since there might be some cells in which the value shown is 0.0 but it is actually greater.

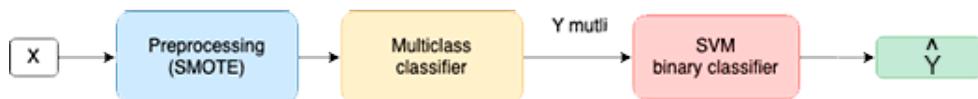
⁹In the standard deviation matrices the results are shown multiplied by 10^2 in order to see the value in the different cells of the matrix

the values obtained are too small to make a good visualization of them so they are just commented.

As mentioned above, we have developed a total of six final implementations. Three of them for a multiclass classification problem that consists of predicting the right label for each sample within the fourteen labels already explained. The other three are basically SVM binary classifiers that take as input the output probabilities of the multiclass classifications. The purpose of these three is to adapt the model to the final objective of creating a system to distinguish between violent and non-violent scenarios. In figure 5.3, two block diagrams are included in order to represent the flow in the multiclass and binary classification approaches.



(a) General model for the three multiclass classification approaches.



(b) Block diagram that shows the concatenation of multiclass classifier and the SVM binary classifier.

Fig. 5.3. Confusion matrices for SVM multiclass classification for train and validation sets

As previously explained in 4.3, as output of the multiclass systems, the soft values were obtained. These are the probabilities of a certain sample to belong to each class, being the highest value the one that usually corresponds to the correct predicted label. So this is exactly what we did: we took the highest probability and consider it the predicted category in order to transform the resulting outputs into hard format. This is the \hat{Y} in the figure.

For the binary approach, we just took the original soft outputs, Y_{multi} , to feed the binary classifier in order to be trained. The model for which the best accuracy value is obtained in the validation set was saved so as to, finally, use it in the testing process.

5.2.1. Implementation 1: SVM classifier for a multiclass classification

We decided to establish a baseline by making a first experiment based on a SVM classifier that is used for a multiclass classification. In this case, we wanted to check the results of the performance of one of the most employed techniques in this field different from NN. It is true that this method is originally designed for binary problems but, as explained in 3.3.2, we took advantage of the multiclass algorithm. To define our model, we used mainly the default parameters. In this work, we have not investigated which type of kernel could better fit our problem, however we found in the literature that the default option of RBF is the most appropriated for real world applications [104].

Below, we can find the results for the different sets. In figure 5.4, the confusion matrices for the train and validation sets are shown. In figure 5.5, the confusion matrix for the evaluation on the test set is included. Finally, in the bar plot of figure 5.6, the accuracy values for the three sets are represented.

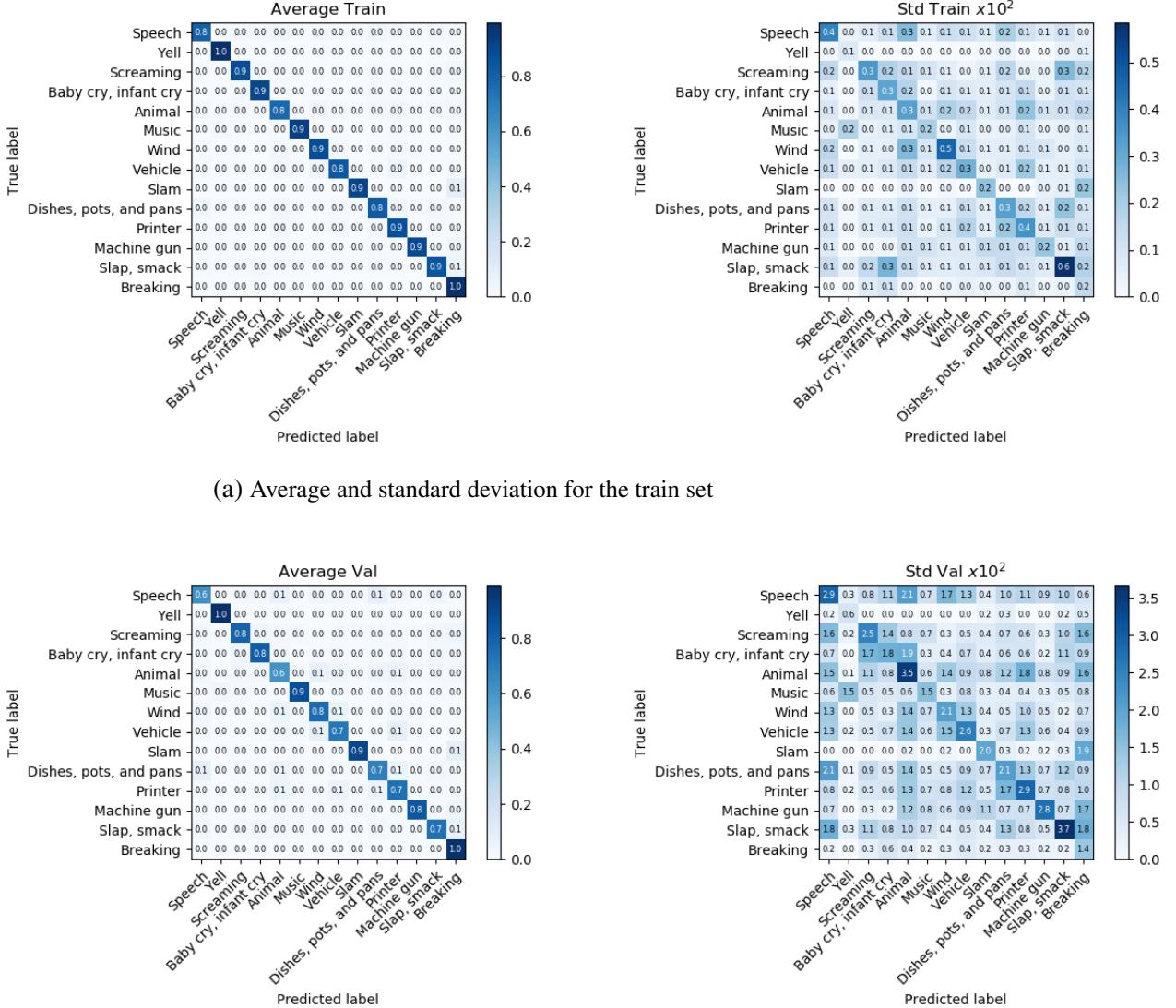


Fig. 5.4. Confusion matrices for SVM multiclass classification for train and validation sets

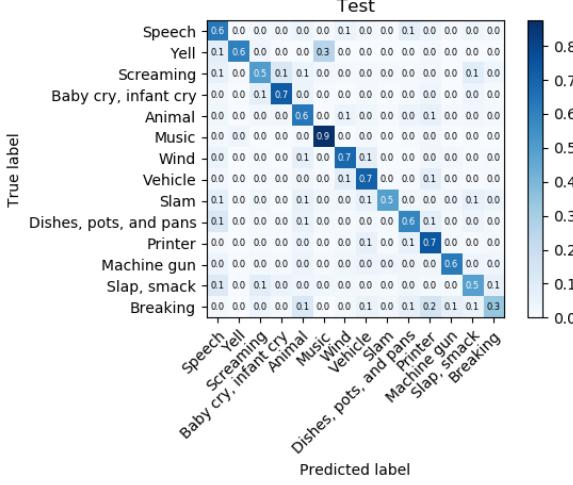


Fig. 5.5. Confusion matrix for the test set for the SVM multiclass classification

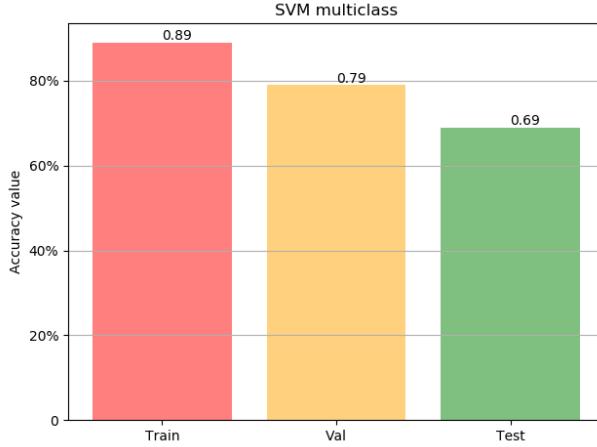


Fig. 5.6. Accuracy values for the three sets for the SVM multiclass classification

The accuracy values are a $89\% \pm 0.087\%$ for the training set, a $79\% \pm 0.56\%$ for the validation and 69% for test. This can also be appreciated in the average confusion matrices. The diagonal for the training and the validation stands out the other cells but they do not present a clear sense of overfitting due to they are not completely uniform. In the test results, we can see greater values in the true positive sections for the most populated classes, mentioned above in section 5.1. This makes sense since the test set has more samples in these categories. We can consider a good result for this case. The accuracy of the final testing process is not much lower than the one obtained for the validation set. It would be a good option to populate more the violent labels, but we did not want to use synthetic data in the final evaluation.

In the table 5.3, a binary grouping considering the violent and non-violent classes is included. This was generated by taking the violent samples predicted as violence, even though the prediction did not match the specific class, it just had to be violent. Same for

non-violent events. Also, the rest of the observations were grouped as predicted violent label for non-violent samples and vice versa.

True/Predicted	Non-violent	Violent
Non-violent	5274	485
Violent	269	870

Table 5.3. ADAPTATION OF THE RESULTS FOR SVM MULTICLASS CLASSIFIER TO BINARY

The accuracy value for this table is 89.07%, which was calculated by using the indicated formula in appendix A. A great difference is notable between the true positive section of the non-violent label with respect to the other parts of the matrix. However, we can consider an accurate result since plenty of the samples belong to true positive regions for all the classes in the confusion matrix for the test set in figure 5.5.

5.2.2. Implementation 2: SVM multiclass classification and SVM binary classification

The SVM multiclass model is then incorporated to the previous stage of the preprocessing. So, it predicts the multiclass probability for each value, what can also be interpreted as it is converting the data to a new feature space before the binary classification.

The results in accuracy and confusion matrices for the different evaluations can be found below. Figure 5.7 includes the average and standard deviation of confusion matrices for train and validation sets. In figure 5.8, the confusion matrix for test is included. Then, in figure 5.9, the accuracy values for the three sets are shown.

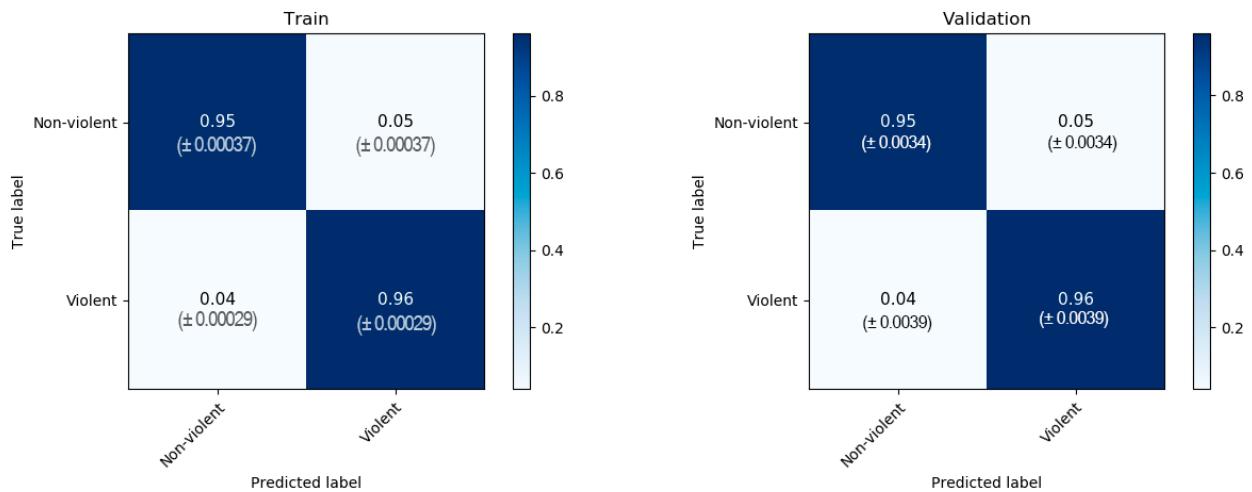


Fig. 5.7. Confusion matrix for the train and validation sets for the SVM multiclass + SVM binary classification

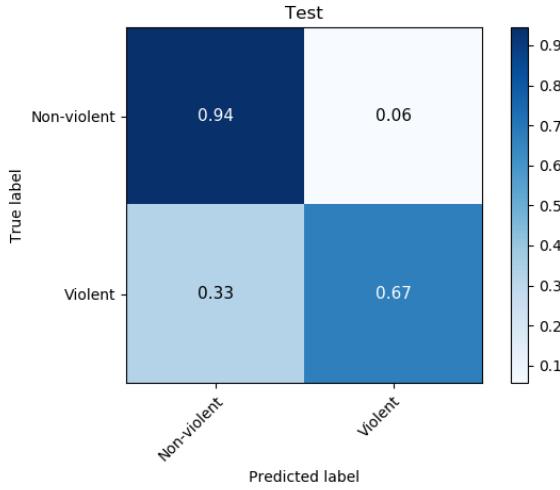


Fig. 5.8. Confusion matrix for the test set for the SVM multiclass + SVM binary classification

In this performing model a more accurate interpretation can be obtained from paying attention to both metrics. In the accuracy bar plot, we can see that the values are $96\% \pm 0.024\%$ for the train set, a $96\% \pm 0.021\%$ for the validation set and a 90% for the test. At first, it can be interpreted that the results are really good. This can be due to the high number of samples for a binary classification. Considering the test value of 90%, we can see that the merit of the good performance is thank to the good classification of the non-violent observations, while the result on differentiating the violent classes is not as good since just the 67% of the labels were correctly predicted. For the train and validation sets, the results are pretty good for both types, being a 96% for both cases. Again, the fact of creating a big amount of the samples artificially for the training and validating parts is given a considerable difference in the results respect to the testing. However, this output also allows to see that with a more balanced test set the performance would be more even in the tree sets, as it happened with the multiclass task explained before, and that those observations which belong to the original embeddings show a satisfying outcome.

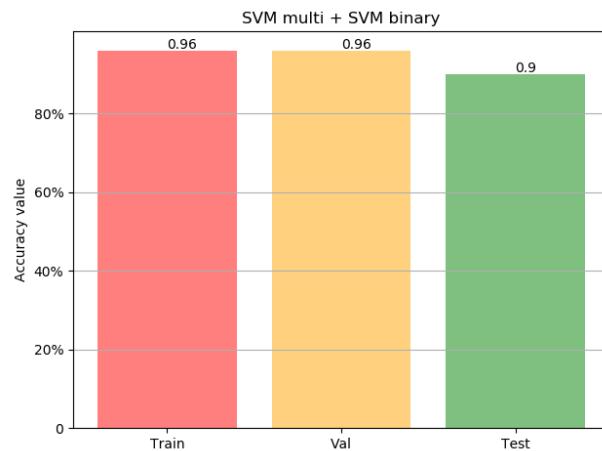


Fig. 5.9. Accuracy values for the three sets for the SVM multiclass + SVM binary classification

5.2.3. Implementation 3: LSTM for multiclass classification

For this implementation we have decided to use a network composed by a total of three layers: two LSTM and a final Fully-Connected with a softmax activation function for the classification task. The first two LSTM layers are set with a drop-out of 0.05 and a recurrent dropout of 0.35. For the first one, the number of units was set to 128, in order not to change the dimensions of its output. In the second layer, it was set to 32 layers so as to reduce the dimensionality before the final prediction in the dense layer. In figure 5.10, an schema of the model is shown.

The function minimized in the process was the categorical cross-entropy, which is a common way to evaluate multiclass classification problems. A more detailed explanation of this function can be found in appendix C. With respect to the training process, a number of 50 epochs were used and also a batch size of 32 samples. Also, a Nadam optimizer was used and a learning rate of $2e^{-3}$. An optimization algorithm is in charge of updating the internal parameters for a better performance and to minimize the loss function. The epochs is a hyperparameter that establishes how many times the learning algorithm pass through the entire train set in the training process. Finally, the batch size is the one that defines the number of observations the algorithm works through before assigns the internal parameters an updated value [105]

Below, the results for the classification are included. In figure 5.11, the average and standard deviation matrices for the train and validation sets are shown. Then, in figure 5.12, the confusion matrix for the testing process is included. In figure 5.13, a bar plot shows the accuracy for every set with their standard deviation values.

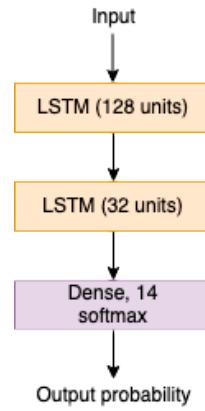
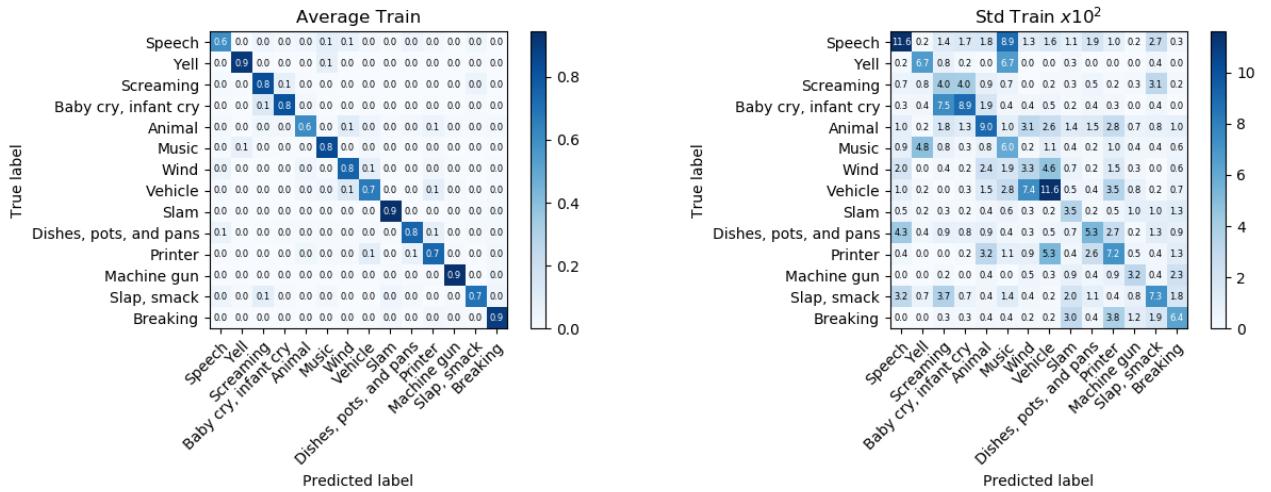
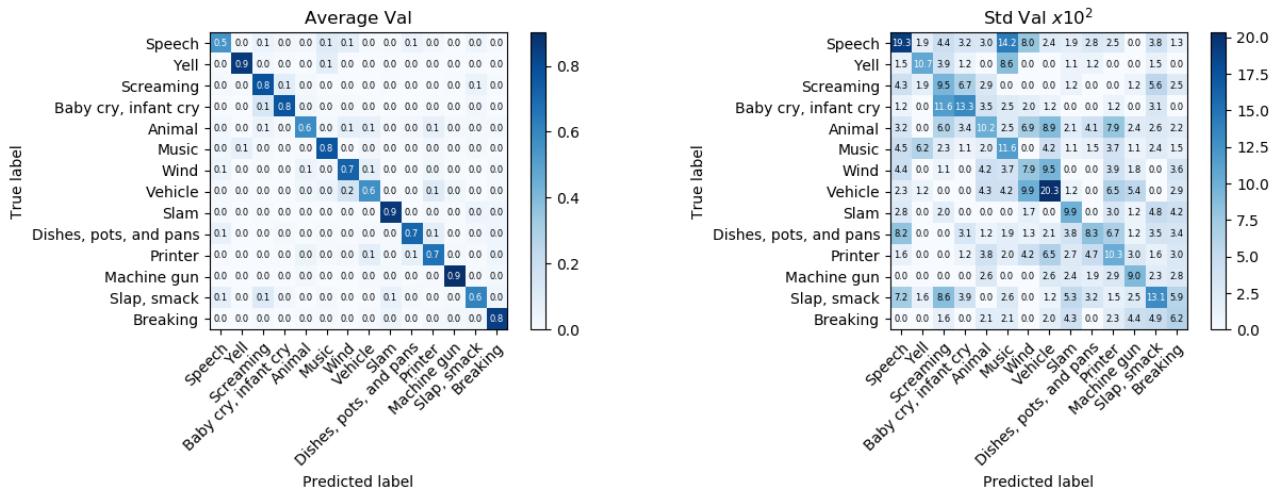


Fig. 5.10. LSTM architecture



(a) Average and standard deviation for the train set



(b) Average and standard deviation for the validation set

Fig. 5.11. Confusion matrices for LSTM multiclass classification for train and validation sets

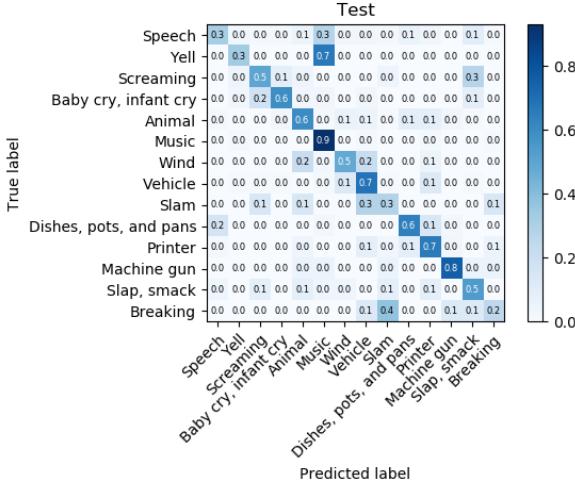


Fig. 5.12. Confusion matrix for the test set for the LSTM multiclass classification

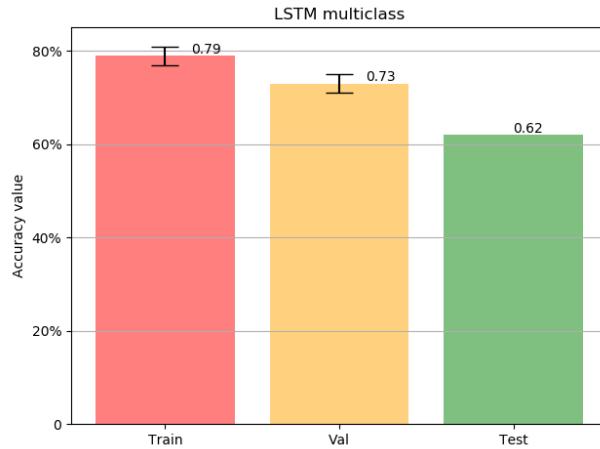


Fig. 5.13. Accuracy values for the three sets for the LSTM multiclass classification

For this model, we obtained an accuracy of $79\% \pm 2.28\%$ for the training set, a $73\% \pm 2.41\%$ for the validation and a 61% for the test. This can be due to the possible confusion regions mainly between classes such as *Screaming* and *Baby cry, infant cry*, and *Wind* and *Vehicle*, whose standard deviation values are the most emphasized. Also, a prominent failing region is the cell *Slam/Vehicle*, showing a not expected behaviour. About the *Yell* class, there is a big misclassification that results in a majority of predicted *Music* labels. The other violent classes have fewer correct predictions but also due to the fewer samples that belong to them in the test subset. We can discard heavy overfitting since the values of accuracy in the training and validation sets with respect to the test one are not too different.

In table 5.4, the adaptation to the binary problem is included. The accuracy for this case is 88.27%. Of course, the true positive region with more samples is again the one that corresponds to non-violent. The observations placed in the false negative cell for the violent class can be a consequence of the misclassification of the samples form *Yell* as

Music. Other cases of failed prediction as *Breaking* for *Slam* do not affect this matrix since the confusion occurs between violent classes.

True/Predicted	Non-violent	Violent
Non-violent	517	62
Violent	20	100

Table 5.4. ADAPTATION OF THE RESULTS FOR LSTM
MULTICLASS CLASSIFIER TO BINARY

5.2.4. Implementation 4: LSTM multiclass + SVM for a final binary classification

For this case, the predicted probabilities obtained from the LSTM multiclass classifier are passed as input to the SVM for the binary classification.

The accuracy value for the different sets are included. Figure 5.14 shows the average and standard deviation of the confusion matrices for the train and validation sets. In figure 5.15, the confusion matrix for the test set is also included, and in figure 5.16, the error bar can be seen with the accuracy values for each of the sets.

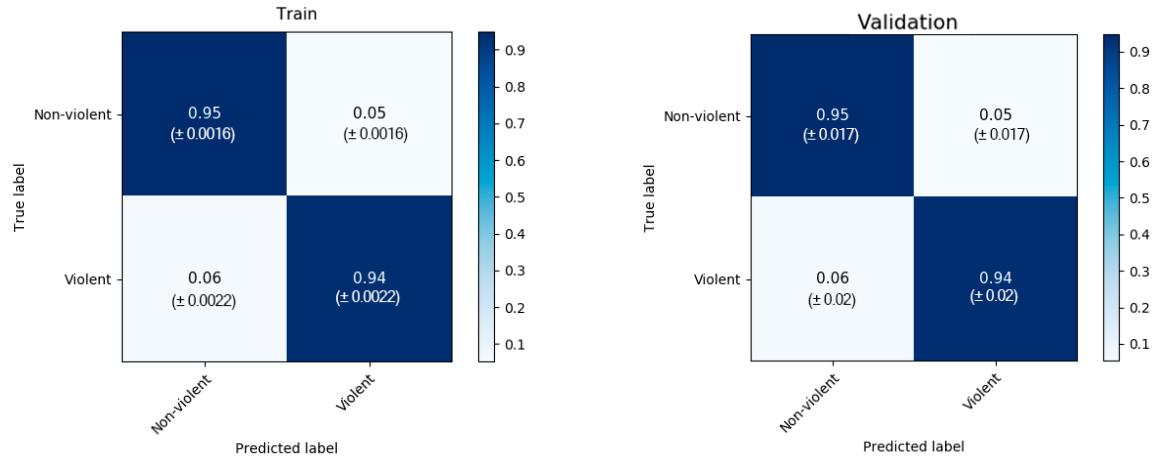


Fig. 5.14. Average and standard deviation for the train and validation set for the LSTM multiclass + SVM binary classification

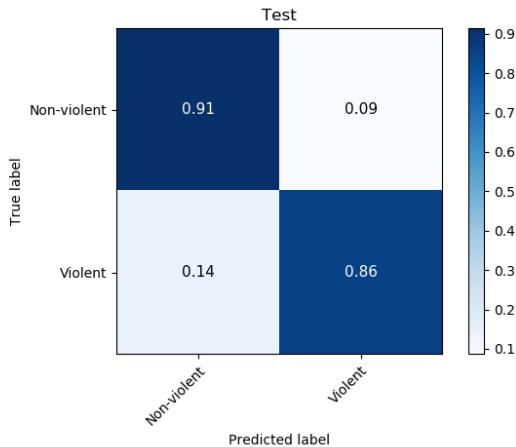


Fig. 5.15. Confusion matrix for the test set for the LSTM multiclass + SVM binary classification

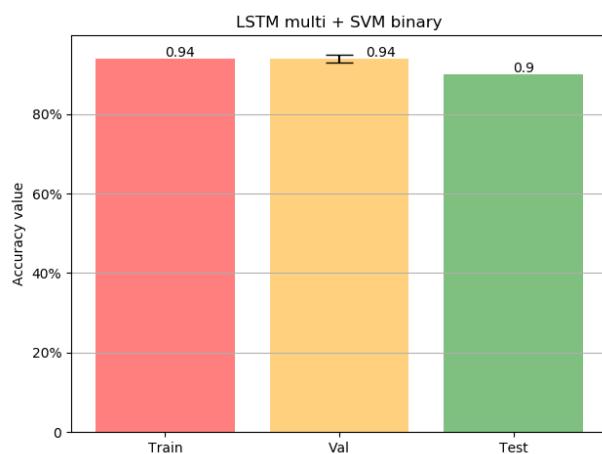


Fig. 5.16. Accuracy values for the three sets for the LSTM multiclass + SVM binary classification

As shown in the bar plot, the values of the accuracy for train and validation are $94\% \pm 0.1\%$ and $94\% \pm 1\%$, respectively. The average matrices are really good as well. The no variation between them could be understood as a symptom of overfitting. However, the accuracy for testing is 90%, which is similar to the train and test results. The interesting point can be read in the confusion matrix for the test set. The classification of the violent classes is almost as good as the one for the non-violent ones, despite of the less amount of samples in these categories.

5.2.5. Implementation 5: CNN for multiclass classification

For this case, we have implemented a CNN model based on the architecture presented in figure 4.3 in subsection 4.2.3 for the small experiment to check the similarity of the embeddings from the *.tfrecord* files and the naturally audio files.

In the configuration of the network, the same hyperparameters from implementation 3 are used for the learning process: a categorical cross-entropy for the loss function, a Nadam optimizer, a number of 50 epochs and a batch size of 32 samples. We have kept the same value for this model since the number of observations is the same and the depth of the network is not very great. A learning rate of $2e^{-3}$ was also used.

In figure 5.17, the average and standard deviation matrices for the train and validation sets are shown. In figure 5.18, the confusion matrix for the test set. The bar plot in 5.19, shows the accuracy for every set with their standard deviation values.

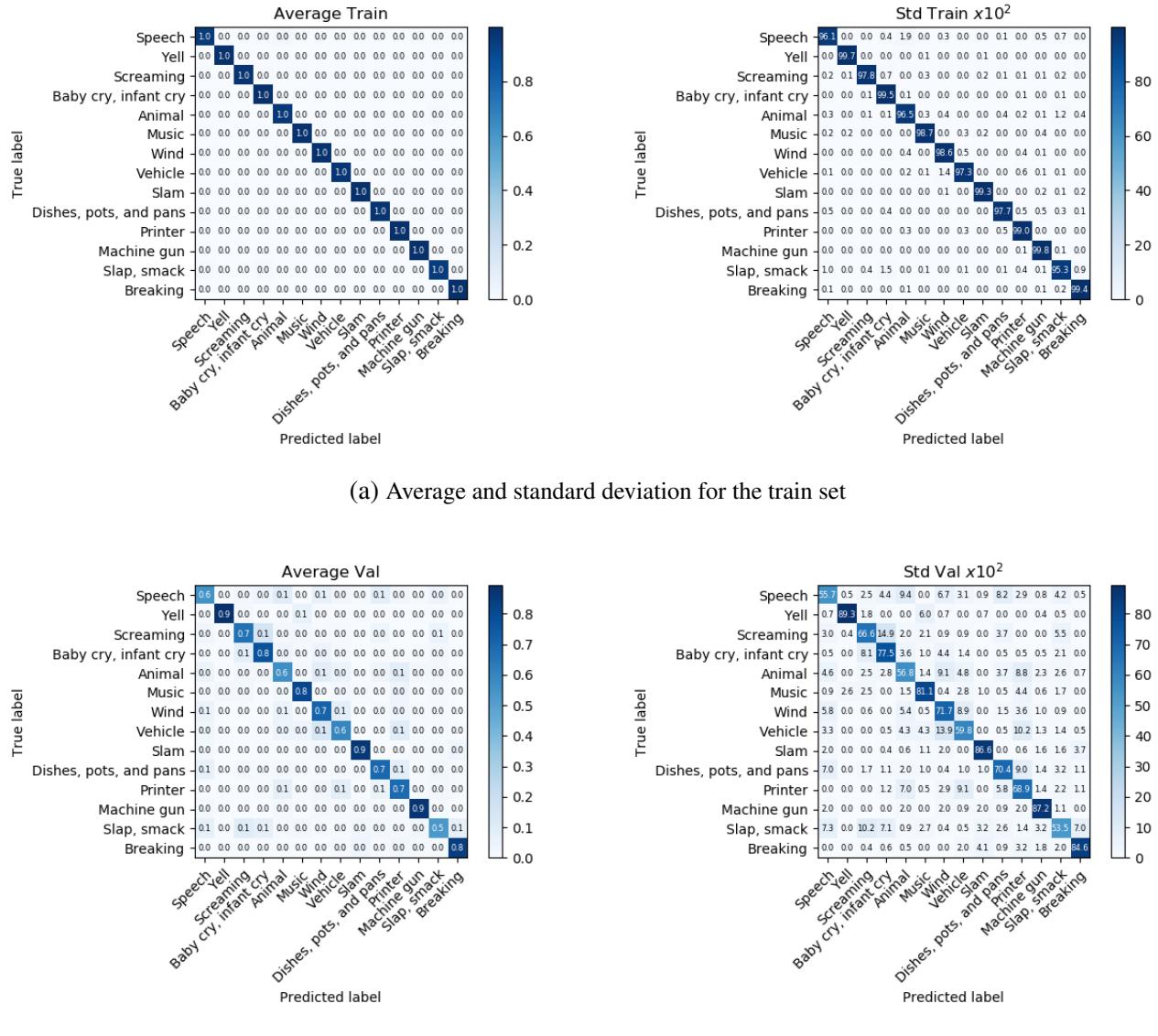


Fig. 5.17. Confusion matrices for CNN multiclass classification for train and validation sets

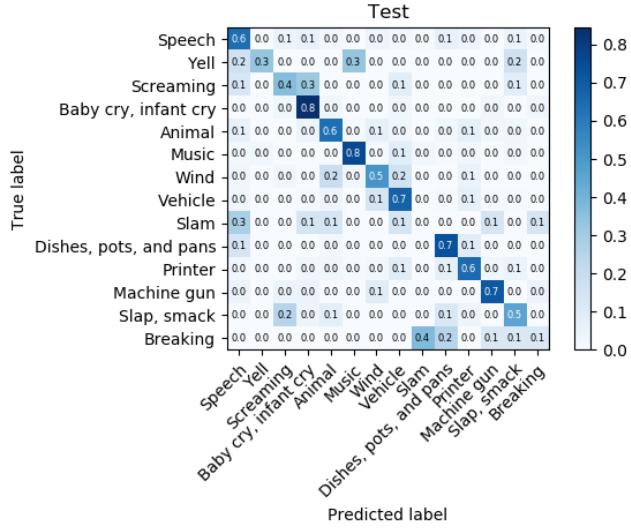


Fig. 5.18. Confusion matrix for the test set for the CNN multiclass classification

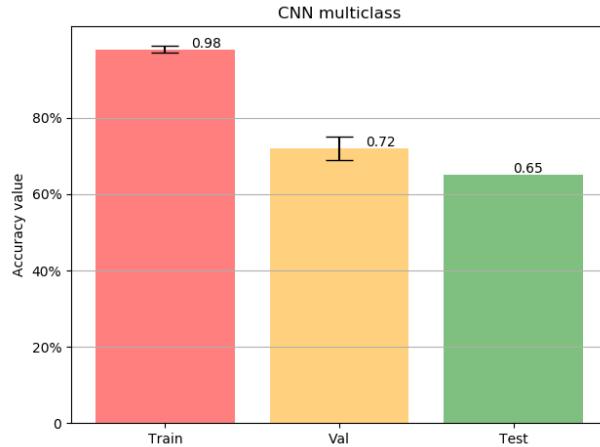


Fig. 5.19. Accuracy values for the three sets for the CNN multiclass classification

The values of the accuracy obtained for this implementation are $98\% \pm 1\%$ for the train set, $72\% \pm 3\%$ for validation and 65% for the test set. In the error bar in figure 5.19, the difference between the red bar for train is considerable meaningful with respect to the other two. Also, the confusion matrix for train shows a very well marked diagonal, while the other two are not that perfect. In the validation set, there are some misclassifications but still the equality among cells is present. Some confusion regions can be found again between *Screaming* and *Baby cry, infant cry* or *Wind* and *Vehicle*. In the test set, we can see that the diagonal does not follow a regular shape, being this one the worst result of the experiments.

As in the previous cases, the table to see the problem from a binary point of view is included in 5.4. For this case the value of the accuracy is 90%. Again this result is kind of delicate. The failed predictions happen between, for example, *Screaming* and *Baby cry, infant cry*, and also between *Breaking* and *Slam*. Errors that this matrix does not show.

True/Predicted	Non-violent	Violent
Non-violent	531	48
Violent	22	98

Table 5.5. ADAPTATION OF THE RESULTS FOR CNN
MULTICLASS CLASSIFIER TO BINARY

5.2.6. Implementation 6: CNN multiclass + SVM for a final binary classification

In this last approach we wanted to use the predictions from the CNN to feed the binary classifier SVM.

Figure 5.20 shows the average and standard deviation of the confusion matrices for the train and validation sets. Figure 5.21 includes the confusion matrix for the test set. Finally, figure 5.22 shows the error bar in which the accuracy values for the three test can be seen.

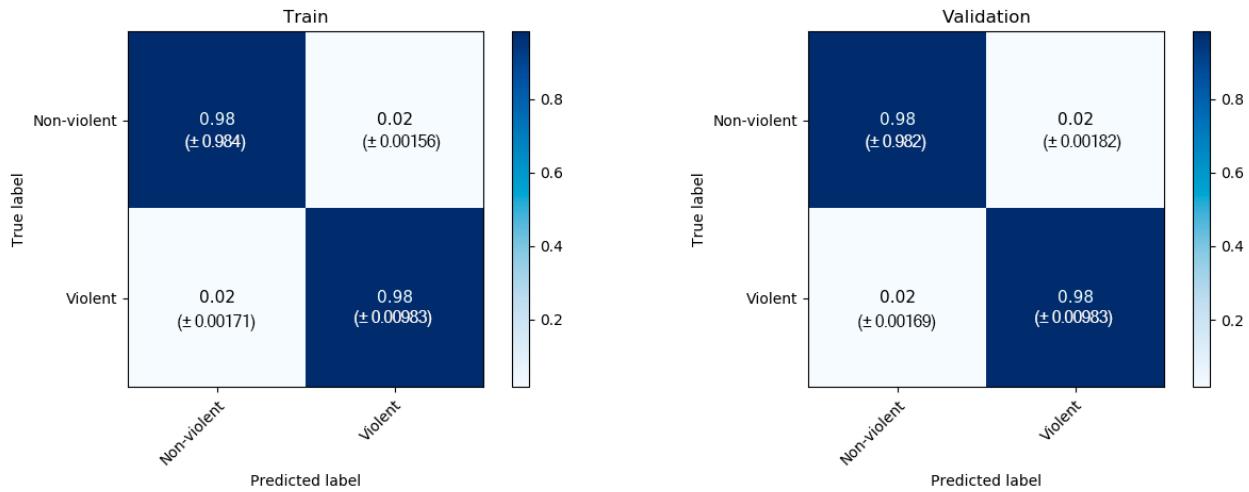


Fig. 5.20. Average and standard deviation for the train and validation set for the CNN multiclass + SVM binary classification

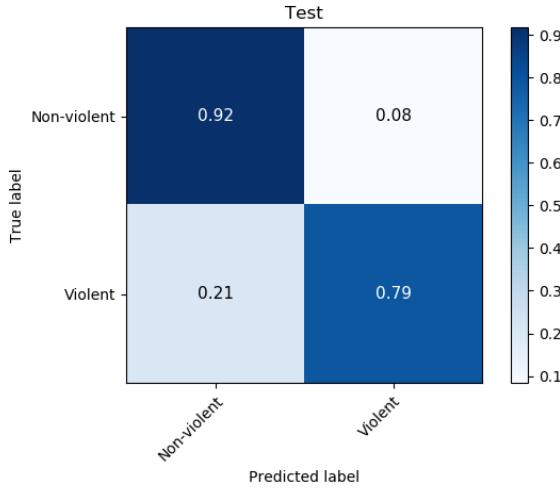


Fig. 5.21. Confusion matrix for the test set for the CNN multiclass + SVM binary classification

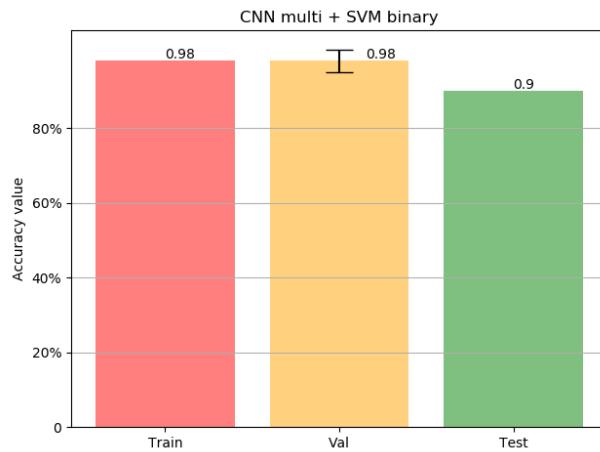


Fig. 5.22. Accuracy values for the three sets for the CNN multiclass + SVM binary classification

The values of the accuracy are $98\% \pm 0.32\%$, $98\% \pm 3\%$ and 90% for train, validation and test set, respectively. Again the results are similar to the ones obtained in the experiment 4. In this case, the classification of the violent classes are good enough as well for the test. However, as seen in experiment 5, there are plenty of misclassifications for the test set but some of them happen between classes of the same kind, so this failures do not have an effect on the binary problem. Also, others have been mixed with categories from the other group, this is would be the reason of the lower output.

5.3. Comparison and results

The idea behind the different implementations was to check the usefulness of the embeddings extracted with VGGish for a violent/non-violent classification problem. We wanted to study the performance of different learning techniques to see how this type of data

worked with them. To do so, we first ran a typical algorithm for classification that does not involve Neural Networks, an SVM. Then, we thought it could be a good idea to try neural models with different core concepts to exploit the nature of our data. The dimensions of our input allowed the use of the CNN, while the sequential character of an audio instance makes the LSTM a possible approach. In this section, we are going to compare the six different results dividing the explanations according to the type classification: multiclass or binary.

For the multiclass problem, the results are shown again in table 5.6.

	Train	Validation	Test
SVM	$92\% \pm 0.087\%$	$86\% \pm 0.56\%$	72%
LSTM	$79\% \pm 2.28\%$	$73\% \pm 2.41\%$	61%
CNN	$98\% \pm 1\%$	$72\% \pm 3\%$	65%

Table 5.6. ACCURACY RESULTS FOR THE THREE DIFFERENT ALGORITHMS AND THE THREE SETS FOR THE MULTICLASS APPROACH

Initially, the best performance for this case is the one done by with SVM. It is expectable that this algorithm achieves a good result due to the way it works. Basically, it expresses the different data samples in a feature space in which the optimal boundary is computed, as explained in subsection 3.3.2. However, this one is not totally comparable with the other two LSTM and CNN, since the number of observations to feed the model is much greater because of the conversions explained in the subsection 5.1. This way, the method has more examples to learn how the different classes are distributed.

For the CNN approach, the results are maybe the least satisfactory. The system presents a clear overfitting due to the big difference of the results for within the three sets. It is learning the data in the training stage and not being able to generalize, so it does not perform a good enough classification for validation and test. Also, considering the errors in figure 5.18, the mismatches are the greatest of the three experiments.

In the LSTM approach, the results are numerically the worst but also the ones that make more sense. The system does not present an overfitting since the accuracies for validation and test are similar to the one for train. Also, most of the errors that can be read in the confusion matrix in figure 5.15 are understandable since the audio data for those classes involved in the confusion regions may look similar. For example, a recording that belongs to *Baby cry, infant cry* is likely to be similar to one from *Screaming*. The same explanation could be used for *Wind* and *Vehicle*, since they usually have sounds that belong to scenes in a rush. Just some cases are more weird as the misclassification of *Yell* for *Music*, or *Screaming* for *Slap, smack*.

Let's look at some violent categories to check their performance.

- The error when predicting *Music* for *Yell* samples appears in all the models in a

more or less clear way. Maybe, the original audio data from which the embeddings were extracted had common similarities. It would not be weird that many samples from *Music* had the presence of shouts and yells.

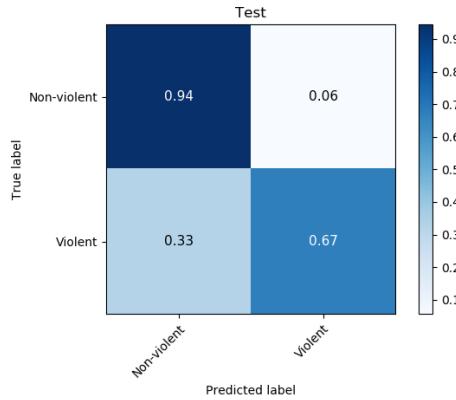
- The confusion region between *Screaming* and *Baby cry, infant cry* comes out in all matrices. This makes sense, since the semantics of these two type of audio events is quite similar.
- The predictions for *Slam* vary in all the models obtaining half of the samples correctly labeled for the SVM approach. In the case of the LSTM the true positive region just counts with 30% of the data but the rest is misclassified among other violent categories. For the CNN, this class is a complete failure being most of it classified as *Speech*.
- In the case of *Machine gun* it has a pretty good performance for all the three models, getting a 80% for LSTM, a 70% for CNN and a 60% in the case of SVM in the true positive regions. The rest of them are shared uniformly along all the other classes.
- In the case of *Slap, smack*, it is predicted in a similar way for all the models as well, achieving a 50% of the samples right classified.
- The class *Breaking* has not showed a really good result for LSTM and CNN being 40% of the samples predicted as *Slam*. However, for SVM a 30% was set as true positives being the rest of the samples uniformly distributed along the other classes, mainly in *Animal* and *Printer*.

The conclusion we can extract from this comparison is that the models which shows much better results are the LSTM and SVM. We can see that the classes have similar results for all the methods being the SVM the one that stands out. In the LSTM, it calls the attention the region between *Wind, Vehicle* and *Slam* and also the amount of samples from *Yell* predicted as *Music*. We could establish that maybe the CNN is not the best option due to the broken diagonal that it presents. Also, a general good performance of the non-violent classes appears in all the matrices. This may be due to the more number of samples that belong to these classes within in the test set. It is worth to mention that the quality of the labels previously described in 3.1.3, can also affect the resulting values. For example, *Yell* presents a quality of the 50% and also *Slam* has a label quality of 20%. This could be taken into account for future implementations.

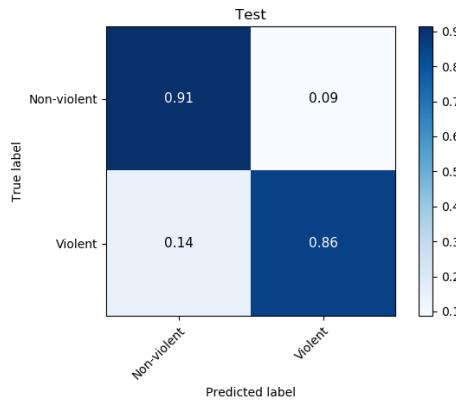
For the binary problem, the results are included in table 5.7. Also the confusion matrices for test are shown again in figure 5.23, since they are really relevant for the final explanation.

	Train	Validation	Test
SVM	$96\% \pm 0.024\%$	$96\% \pm 0.021\%$	90%
LSTM	$94\% \pm 0.1\%$	$95\% \pm 1\%$	90%
CNN	$98\% \pm 0.32\%$	$98\% \pm 3\%$	90%

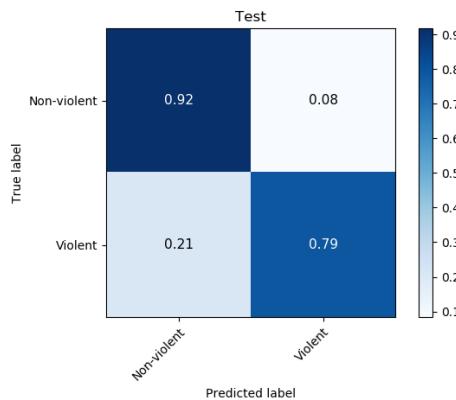
Table 5.7. ACCURACY RESULTS FOR THE THREE DIFFERENT ALGORITHMS AND THE THREE SETS FOR THE BINARY APPROACH



(a) SVM multiclass + SVM binary



(b) LSTM multiclass + SVM binary



(c) CNN multiclass + SVM binary

Fig. 5.23. Test confusion matrices

The actual solution of the problem would be this binary classification that allows us to distinguish between non-violent and violent events. The results apparently are really good if we look at the accuracy values, but the interpretability is much more clear when reading the confusion matrices.

In order to interpret the result for the SVM, we can read the confusion matrix in figure 5.5 to check how the multiclass classification performed and related to the output obtained from the binary classifier. For this case, certain classes have can be predicted within the range of false negatives as non-violent categories, such as *Yell* confused with *Music*. Also, a couple of categories are predicted as *Speech* and some others as *Animal*. If we interpret this errors in a binary context, we can say that the system is taking advantage of the higher a priory of non-violent classes without learning the true distinctions. This fact is reflected in the only 67% of true positives for the violent label in figure 5.23 (a).

In the CNN approach, the results are numerically smaller but the classification of the violent instances is much better. Almost a 80% of these samples were correctly classified. Again, if we look at figure 5.18, some classes are confused within the two different binary labels. A 30% of instances from *Slam* have been wrongly identified as *Speech*. The same happens with *Yell* and *Music* as in the case of SVM. So, these can affect the output of this binary approach and the result of the test confusion matrix. There is also a big mistake in predicting *Slam* for samples that belong to *Breaking*, but this fact does not involve errors in the binary classification.

Finally, the LSTM is the method for which a better result is obtained in this problem. In the confusion matrix in figure 5.12, as explained above, we can see some misclassifications but they happen between data of the same type, except between *Yell* and *Music*. In a real application this would be a more reasonable result, in order to detect a violent situation without depending too much on the exact event occurred just considering if it is violent or not. This is why the percentage of true positive samples for the violent class is the best for this model with a 86%. We can say that this is the best approach from all of our experiments.

6. CONCLUSION

A couple of conclusions can be extracted from the whole development of the work related to the different objectives that were initially set. Within the realization of the steps that were considered necessary to achieve the solution of the final task, several outcomes that should be taken into account were achieved.

During the first objective of finding a proper database that matched the problem, we came up with several solutions. Finally, Audio Set was the best choice due to its characteristics of containing a massive amount of videos in different conditions. Initially, we thought that the lack of data would not be a problem and that we could extract all the necessary information from all these videos. However, some limitations were discovered during the whole investigation process. The unbalanced nature of the dataset, the limitation that we imposed of just using single label data and the quality of the different labels are some examples of the problems that we found when adapting Audio Set to our problem. Also, the own quality of the videos could be a factor that may affect the resulting embedding. Recording a more suitable dataset could be a good point to improve since this is one of the work packages of the project EMPATIA.

However, we did not use this database only as a data resource. The same developers that made Audio Set, released a deep neural model called VGGish which certainly eased the labour of extracting features. We investigated this network and checked the behaviour of the feature embeddings and concluded that they could actually be used to perform classification tasks within this research area. Maybe, we cannot consider this transformation as a feature space in which the violent component plays a really significant role when finding a separability between violent and non-violent events, but the feature extractor we used is a good option that should be considered when implementing future works.

Regarding the implementation of a violent event classifier from a gender-based violence point of view, it is clear that the definition of how these events must look is fundamental. The system designed for that purpose is just a simple method to select violent classes in a more efficient way whose usefulness validation is out of the scope of this work. Considering the well structured ontology that Audio Set is provided us with, we thought that a good way of saving time would be to offer the possibility of selecting violent categories without checking all the labels one by one by navigating the ontology. However, the idea is mainly the same that a normal application or interacting system should follow to accomplish this task, maybe with other type of classes or more specific ones.

In order to implement a learning method to classify violent events, we decided to use a cascade of a multiclass classification and a binary one. This solution was not the one that we intend at the beginning. One of the initial approaches was to perform a normal binary model with data tagged as *violent* or *non-violent*. However, we thought it could be a good idea to focus on obtaining the multiclass results and then adapt them to the

binary situation and see the outcome. With this purpose, we compared different methods and results, being the SVM model the one with a best performance on the multiclass task and the LSTM the one with most accurate results for the binary case. However, a final conclusion of which one is better cannot be extracted. In the multiclass classification, the output obtained generally made sense if considering the semantics of the label even though the true positive regions were not the most populated. In the binary classes, the results can be considered accurate but the number of classes is much higher and also most of the corrected predictions belonged to non-violent events, since these are the most populated classes in the test set. The CNN model did not stand out in any of the tasks, but it achieved pretty good results in the binary classification. The architecture of the neural models is an aspect that also should be taken into account, since this work did not go deep enough in that regard.

7. FUTURE WORK

As explained in the section 1.2, this work can be defined as the very first approach of one of the parts of a much larger and ambitious project. So, it is easy to understand the future work is very extensive in order to fulfill the final requirements. However, we propose a couple of details and considerations that we thought a good idea when trying to implement our approach.

One of the most important aspects when working on machine and deep learning problems is the quality and amount of data. Finding these resources is actually a really tough and long task. This is one of the first points that should be taken into account in order to go further in this investigation line. Researchers from the EMPATIA-TC project are currently working to solve this problem by trying to develop a database that collects data related to stimuli, labels for emotions and reactions of victims [5] in situations of violence against women. However, within the scope of this work, since the violence databases that are available do not focus on this task or are little populated, it would be interested to make a public dataset that allow to fully adapt the problem to the gender-based violence, including a wide range of sounds that may really happen in this kind of situations. This could be done in a much less complicated way without recording any victims, for example, by collecting audio data that strictly belong to this field from public websites and subject them to a proper labeling process.

More work on trying to establish a subjective definition of violence that could be implemented is also possible. For our case, considering the way we chose, explained in subsection 2.3.2, a better implementation could be designed to ease the contribution of the victim to the whole learning process so it is not so annoying to be selecting all the categories, or, maybe, an automatic system in which the victim does not need to directly participate.

Related to the feature extraction process, more studies could be done in order to finally select a feature space that totally adapts to the problem and makes the separability between violent and non-violent events available. If we consider the embeddings extracted with VGGish, a more immediate research line could be to try different values of the parameters related to the final shape of the embedding and overlapping between rows, as mentioned in subsection 3.1.4, so as to see the different results when extracting embeddings from audio data.

A good study to dive into is to improve the system used in the learning process by trying to apply new methods and algorithms or just taking more advantage of the ones already used, explained in subsection 3.3 and chapter 5. For example, adapting the CNN architecture model to the available data or improving the performance of the LSTM. In this last aspect, we did try to study a little the possibility of implementing an attention-based approach mixed with LSTM, following the idea proposed in [106] for sentiment

analysis, but it was going to become very complex and out of our initial scope.

Finally, the implementation of a model for multilabel¹⁰ classification. This is actually a very important topic that could increase the scope of the problem, also paying attention to its real world application. We just worked with single label samples from the dataset, which supposed a big limitation on the process of class selection. A bit of research was also done in order to make an attempt on this task by trying to use Classifier-Chains (CC), which treats the multilabel problem as a combination of several single label classifications [107], but we did not get any further.

¹⁰Multilabel classification differs from multiclass classification in the way the samples are labeled. The former treats samples that belong to more than one class at the same time, and the latter, the one implemented in this project, works with samples that belong to just one class.

8. PLANNING AND ECONOMIC RESOURCES

8.1. Planning and time organization

This whole project has been initially assigned in November 2018. Some research was done during the first half of 2019, but it was not fully addressed until September 2019. In order to explain the development of the tasks, these are listed below:

- Resources research: tasks related to the process of finding a database and a method to extract features. Also, the understanding of the whole feature extractor we use in the work can be included here.
- Previous works investigation: research about the current situation of the audio classification and how to implement a system with the obtained features.
- First attempts and experimentation: failed experiments and evolution of the designed system solution.
- Implementation of the final model: building the final system proposed and work to obtain the results.
- Assistance to meetings: assistance to meetings about the EMPATIA project
- Writing the report

In table 8.1 a relation between the estimated time per task and the final amount of hours invested is included. It must be taken into account that in this time we are considering as well the time between the first period of the work defined above.

Task	Estimation (h)	Actual time (h)
Resources research	70	100
Previous works investigation	100	120
First attempts and experimentations	100	200
Implementation of the final model	160	160
Assistance to meetings	0	4
Writing the report	320	320
TOTAL	750	904

Table 8.1. RELATION OF ESTIMATED TIME AND ACTUAL TIME
INVESTED PER TASK IN HOURS

8.2. Economic resources

There are two main types of cost in to the development of this task. One of them is related to the material use and the other one to the recruitment. With respect to the first one, the only material that has been used is a computer *MacBook Pro (Retina, 13 inches, End of 2013)*. The equivalent computer nowadays has a total price of 1.203,84 euros. Consider that a normal computer has a lifetime of 5 years of duration, that a month is composed by 160 hours of work and the number of total hours of work have been 904, we can deduce that the proportional price of the device is 113,36 euros without taxes.

As for the staff costs, the people that have participated in this task is the student, as a researcher, and the professor, as a project manager. Also, some participation from the people in the TSC department can be included. Considering an average salary of a researcher in Spain of 1.286 euros per month, the salary for the student in case of being an official researcher would be 7.225,90 euros. Also, for the two participants of the department, considering they spent a total of 3 hours each, a salary of 24,1125 euros would correspond to the help of the researches [108].

About the professor, considering that the period in which the project took place is approximately a year and that meetings of half an hour with the student took place averagely one time per week, if we also discard the two months of summer, an amount of 60 hours in meetings were invested by the professor. Also considering an extra time of 20 hours in the last part when working in the report, this gives a total of 80. If the average salary per month of a project manager in Spain is 2.217 euros, then a salary of 1.108,50 euros would correspond to the professor [109].

In table 8.2 the different expenses and salaries are included and also a final budget.

Resource or worker	Cost (euros)
MacBook Pro	113,36
Assistant researcher	24,1125
Assistant researcher	24,1125
Main researcher	7.225,90
Professor	1.108,50
TOTAL	8.495,99

Table 8.2. RELATION OF ESTIMATED TIME AND ACTUAL TIME
INVESTED PER TASK IN HOURS

ACRONYMS

.csv Comma-separated values. 20, 21, 40, 41

.wav Waveform Audio File Format. 40

AED Acoustic Event Detection. 7, 9, 16

AED/C Acoustic Event Detection and Classification. iii, vii, 4

ANN Artificial Neural Network. vii, 9, 17, 21, 22, 26, 33

ASC Acoustic Scene Classification. iii, vii, 4, 7, 8, 13

AVD Acoustic Violent Detection. iii

BAA Broad Agency Announcement. 11

BOF Bag-of-frames. 7, 8

CASA Computational Acoustic Scenes Analysis. 4

CC Classifier-Chains. 75

CHIL Computers in the Human Interaction Loop. 16

CNN Convolutional Neural Network. iv, vii, viii, 9, 10, 17, 21, 22, 23, 26, 27, 30, 31, 42, 48, 49, 50, 53, 63, 66, 68, 69, 70, 71, 73, 74

ConvNet Convolutional Network. 22, 23, 24, 25, 27, 28

CV Computer Vision. 4, 9, 27

D Dimensions. 23, 24, 25, 48

DARPA Defense Advanced Research Projects Agency. 11

DFT Discrete Fourier Transform.

DNN Deep Neural Network. 10, 13, 22

FC Fully-Connected. 9, 21, 22, 23, 26, 27, 28, 29, 43, 59

FFT Fast Fourier Transform. 6,

FT Fourier Transform.

GMM Gaussian Mixture Models. 7, 9

- HMM** Hidden Markov Models. 8, 9
- ILSVRC** ImageNet Large Scale Visual Recognition Challenge. 27, 28
- IPTO** Information Processing Technology Office. 11
- IPV** Intimate Partner Violence. iii, 1, 14
- k-Fold** k-Fold Cross Validation. 53,
- L³** Look, Listen and Learn. 10
- LLD** low-level descriptors. iii, 5, 6, 7, 8
- LRN** Local Response Normalisation. 23, 28
- LSTM** Long Short Term Memory. iv, viii, 9, 10, 33, 34, 48, 49, 50, 53, 59, 62, 68, 69, 70, 71, 73, 74
- MED** Multimedia Event Detection. iii, 7
- MFCC** Mel-frequency Cepstrum Coefficients. 6, 7, 9,
- MID** Machine ID. 19, 20
- MIR** Music Information Retrieval. iii
- MLP** multi-layer perceptron. 22
- Nadam** Nesterov-accelerated Adaptive Moment Estimation. 59, 64
- NN** Neural Network. iv, 10, 39, 44, 54, 68
- PCA** Principal Components Analysis. 31, 44
- RBF** Radial Basis function. 32, 54
- ReLU** Rectified Linear Unit. 24, 28, 43
- RGB** red, green and blue. 23, 27
- RNN** Recurrent Neural Network. viii, 9, 33, 34, 35, 36, 48
- SC** Spectral Centroid. 6
- SDG** Sustainable Development Goals. 3
- SF** Spectral Flux. 6, 7
- SMOTE** Synthetic Minority Over-sampling Technique. vii, 11, 31, 48, 51

- SNE** Stochastic Neighbor Embedding. 45
- STFT** Short-Time Fourier Transform. 6, 30,
- SVC** C-Support Vector Classification. 32
- SVM** Support Vector Machine. iv, vii, viii, x, 8, 9, 32, 33, 48, 49, 50, 51, 53, 54, 57, 62, 66, 68, 69, 71, 73
- t-SNE** t-distributed Stochastic Neighbor Embedding. 42, 45
- tanh** Hyperbolic tangent. 35, 36
- TF** TensorFlow. 20, 43
- UC3M** Universidad Carlos III de Madrid. 14
- UN** United Nations. 3
- URL** Uniform Resource Locator. 40
- VED** Violent Event Detection. iv, vii, 12
- VGG** Visual Geometry Group. iv, vii, viii, 10, 12, 17, 21, 27, 29, 30, 38, 39, 41, 67, 72, 74
- ZCR** Zero Crossing Rate. 6

BIBLIOGRAPHY

- [1] L. Heise, M. Ellsberg, and M. Gottemoeller, *Ending violence against women*. 1999. doi: [10.4324/9780429269516-5](https://doi.org/10.4324/9780429269516-5).
- [2] C. Watts and C. Zimmerman, *Violence against women: Global scope and magnitude*, 2002. doi: [10.1016/S0140-6736\(02\)08221-1](https://doi.org/10.1016/S0140-6736(02)08221-1).
- [3] P. Blanco, C. Ruiz-Jarabo, L. García de Vinuesa, and M. Martín-García, “La violencia de pareja y la salud de las mujeres”, *Gaceta Sanitaria*, 2004. doi: [10.1157/13062524](https://doi.org/10.1157/13062524).
- [4] World Health Organization, “Understanding and addressing violence against women: Intimate partner violence”, *Understanding and addressing violence against women*, 2012.
- [5] UC3M4Safety Team, *protEcción integral de las víctimas de violencia de género Mediante comPutación AfecTIva multimodAl*, 2018.
- [6] L. Baldwin, *Aggravated Sexual Assault Charges and Penalties*.
- [7] Goetzparterns, *Privacy Policy YouTube*.
- [8] United Nations, *Goal 5: Achieve gender equality and empower all women and girls*.
- [9] D. Barchiesi, D. D. Giannoulis, D. Stowell, and M. D. Plumbley, “Acoustic Scene Classification: Classifying environments from the sounds they produce”, *IEEE Signal Processing Magazine*, 2015. doi: [10.1109/MSP.2014.2326181](https://doi.org/10.1109/MSP.2014.2326181).
- [10] D. Dubois, C. Guastavino, and M. Raimbault, “A cognitive approach to urban soundscapes: Using verbal data to access everyday life auditory categories”, *Acta Acustica united with Acustica*, 2006.
- [11] D. Wang and G. J. Brown, *Computational auditory scene analysis: Principles, algorithms, and applications*. 2006. doi: [10.1109/9780470043387](https://doi.org/10.1109/9780470043387).
- [12] A. J. Eronen *et al.*, “Audio-based context recognition”, in *IEEE Transactions on Audio, Speech and Language Processing*, 2006. doi: [10.1109/TSA.2005.854103](https://doi.org/10.1109/TSA.2005.854103).
- [13] M. Bahoura, “Pattern recognition methods applied to respiratory sounds classification into normal and wheeze classes”, *Computers in Biology and Medicine*, 2009. doi: [10.1016/j.combiom.2009.06.011](https://doi.org/10.1016/j.combiom.2009.06.011).
- [14] D. Van Nort, P. Oliveros, and J. Braasch, “Electro/acoustic improvisation and deeply listening machines”, *Journal of New Music Research*, vol. 42, no. 4, pp. 303–324, 2013.

- [15] A. Temko *et al.*, “Acoustic Event Detection and Classification”, in *Computers in the Human Interaction Loop*, 2009, ch. Part II, 7. doi: [10.1007/978-1-84882-054-8_7](https://doi.org/10.1007/978-1-84882-054-8_7).
- [16] A. Temko *et al.*, “CLEAR evaluation of acoustic event detection and classification systems”, in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2007. doi: [10.1007/978-3-540-69568-4_29](https://doi.org/10.1007/978-3-540-69568-4_29).
- [17] E. S. Sazonov *et al.*, “Automatic detection of swallowing events by acoustical means for applications of monitoring of ingestive behavior”, *IEEE Transactions on Biomedical Engineering*, 2010. doi: [10.1109/TBME.2009.2033037](https://doi.org/10.1109/TBME.2009.2033037).
- [18] I. Potamitis, S. Ntalampiras, O. Jahn, and K. Riede, “Automatic bird sound detection in long real-field recordings: Applications and tools”, *Applied Acoustics*, 2014. doi: [10.1016/j.apacoust.2014.01.001](https://doi.org/10.1016/j.apacoust.2014.01.001).
- [19] Y. Wang, L. Neves, and F. Metze, “Audio-based multimedia event detection using deep recurrent neural networks”, in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 2016. doi: [10.1109/ICASSP.2016.7472176](https://doi.org/10.1109/ICASSP.2016.7472176).
- [20] T. Giannakopoulos and A. Pikrakis, *Introduction to Audio Analysis: A MATLAB Approach*. 2014. doi: [10.1016/C2012-0-03524-7](https://doi.org/10.1016/C2012-0-03524-7).
- [21] X. Amatriain, “An Object-Oriented Metamodel for Digital Signal Processing with a focus on Audio and Music”, 2004.
- [22] D. Marr, “Vision: a computational investigation into the human representation and processing of visual information.”, *Vision: a computational investigation into the human representation and processing of visual information.*, 1982. doi: [10.1016/0022-2496\(83\)90030-5](https://doi.org/10.1016/0022-2496(83)90030-5).
- [23] T. Giannakopoulos, D. Kosmopoulos, A. Aristidou, and S. Theodoridis, “Violence content classification using audio features”, in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2006. doi: [10.1007/11752912_55](https://doi.org/10.1007/11752912_55).
- [24] J. García-Gómez, M. Bautista-Durán, R. Gil-Pita, I. Mohino-Herranz, and M. Rosa-Zurera, “Violence detection in real environments for smart cities”, in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2016. doi: [10.1007/978-3-319-48799-1_52](https://doi.org/10.1007/978-3-319-48799-1_52).
- [25] O. BÜYÜK and M. L. ARSLAN, “Combination of Long-Term and Short-Term Features for Age Identification from Voice”, 2018.
- [26] D. Stowell, D. Giannoulis, E. Benetos, M. Lagrange, and M. D. Plumbley, “Detection and Classification of Acoustic Scenes and Events”, *IEEE Transactions on Multimedia*, 2015. doi: [10.1109/TMM.2015.2428998](https://doi.org/10.1109/TMM.2015.2428998).

- [27] J. R. Uijlings, A. W. Smeulders, and R. J. Scha, “Real-time Bag of Words, approximately”, in *CIVR 2009 - Proceedings of the ACM International Conference on Image and Video Retrieval*, 2009. doi: [10.1145/1646396.1646405](https://doi.org/10.1145/1646396.1646405).
- [28] J.-J. Aucouturier, B. Defreville, and F. Pachet, “The bag-of-frames approach to audio pattern recognition: A sufficient model for urban soundscapes but not for polyphonic music”, *The Journal of the Acoustical Society of America*, 2007. doi: [10.1121/1.2750160](https://doi.org/10.1121/1.2750160).
- [29] J.-J. Aucouturier and B. Defreville, “Judging the similarity of soundscapes does not require categorization: Evidence from spliced stimuli”, *The Journal of the Acoustical Society of America*, 2009. doi: [10.1121/1.3083232](https://doi.org/10.1121/1.3083232).
- [30] A. Mesaros, T. Heittola, A. Eronen, and T. Virtanen, “Acoustic event detection in real life recordings”, in *European Signal Processing Conference*, 2010.
- [31] T. Grill, “Constructing high-level perceptual audio descriptors for textural sounds”, in *Proceedings of the 9th Sound and Music Computing Conference, SMC 2012*, 2012.
- [32] F. Pachet and P. Roy, “Analytical features: A knowledge-based approach to audio feature generation”, *Eurasip Journal on Audio, Speech, and Music Processing*, 2009. doi: [10.1155/2009/153017](https://doi.org/10.1155/2009/153017).
- [33] D. Bogdanov, J. Serrà, N. Wack, P. Herrera, and X. Serra, “Unifying low-level and high-level music similarity measures”, *IEEE Transactions on Multimedia*, 2011. doi: [10.1109/TMM.2011.2125784](https://doi.org/10.1109/TMM.2011.2125784).
- [34] H. Jiang, J. Bai, S. Zhang, and B. Xu, “SVM-based audio scene classification”, in *Proceedings of 2005 IEEE International Conference on Natural Language Processing and Knowledge Engineering, IEEE NLP-KE'05*, 2005. doi: [10.1109/NLPKE.2005.1598721](https://doi.org/10.1109/NLPKE.2005.1598721).
- [35] J. T. Geiger, B. Schuller, and G. Rigoll, “Large-scale audio feature extraction and SVM for acoustic scene classification”, in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2013. doi: [10.1109/WASPAAC.2013.6701857](https://doi.org/10.1109/WASPAAC.2013.6701857).
- [36] Z. Fu, G. Lu, K. M. Ting, and D. Zhang, “A survey of audio-based music classification and annotation”, *IEEE Transactions on Multimedia*, 2011. doi: [10.1109/TMM.2010.2098858](https://doi.org/10.1109/TMM.2010.2098858).
- [37] S. J. Kwon, *Artificial neural networks*. 2011. doi: [10.4324/9781315154282-3](https://doi.org/10.4324/9781315154282-3).
- [38] X. Zhuang, X. Zhou, M. A. Hasegawa-Johnson, and T. S. Huang, “Real-world acoustic event detection”, *Pattern Recognition Letters*, 2010. doi: [10.1016/j.patrec.2010.02.005](https://doi.org/10.1016/j.patrec.2010.02.005).
- [39] Z. Ren, Q. Kong, K. Qian, M. D. Plumley, and B. W. Schuller, “ATTENTION-BASED CONVOLUTIONAL NEURAL NETWORKS FOR ACOUSTIC SCENE CLASSIFICATION”, 2018.

- [40] J. Cramer, H. H. Wu, J. Salamon, and J. P. Bello, “Look, Listen, and Learn More: Design Choices for Deep Audio Embeddings”, in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 2019. doi: [10.1109/ICASSP.2019.8682475](https://doi.org/10.1109/ICASSP.2019.8682475).
- [41] Y. Aytar, C. Vondrick, and A. Torralba, “SoundNet: Learning sound representations from unlabeled video”, in *Advances in Neural Information Processing Systems*, 2016. arXiv: [1610.09001](https://arxiv.org/abs/1610.09001).
- [42] J. Salamon, D. MacConnell, M. Cartwright, P. Li, and J. P. Bello, “Scaper: A library for soundscape synthesis and augmentation”, in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2017. doi: [10.1109/WASPAA.2017.8170052](https://doi.org/10.1109/WASPAA.2017.8170052).
- [43] H. K. Jabbar and R. Z. Khan, “Methods to Avoid Over-Fitting and Under-Fitting in Supervised Machine Learning (Comparative Study)”, 2015. doi: [10.3850/978-981-09-5247-1_017](https://doi.org/10.3850/978-981-09-5247-1_017).
- [44] T. Ko, V. Peddinti, D. Povey, and S. Khudanpur, “Audio augmentation for speech recognition”, in *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, 2015.
- [45] J. Salamon and J. P. Bello, “Deep Convolutional Neural Networks and Data Augmentation for Environmental Sound Classification”, *IEEE Signal Processing Letters*, 2017. doi: [10.1109/LSP.2017.2657381](https://doi.org/10.1109/LSP.2017.2657381). arXiv: [1608.04363](https://arxiv.org/abs/1608.04363).
- [46] S. Wei *et al.*, “SAMPLE MIXED-BASED DATA AUGMENTATION FOR DOMESTIC AUDIO TAGGING”, *Detection and Classification of Acoustic Scenes and Events 2018*, 2018.
- [47] S. J. Pan and Q. Yang, *A survey on transfer learning*, 2010. doi: [10.1109/TKDE.2009.191](https://doi.org/10.1109/TKDE.2009.191).
- [48] D. Sarkar, “A Comprehensive Hands-on Guide to Transfer Learning with Real-World Applications in Deep Learning”, 2018.
- [49] S. Ruder, “Transfer Learning - Machine Learning’s Next Frontier”, 2017.
- [50] E. G. Krug, J. A. Mercy, L. L. Dahlberg, and A. B. Zwi, “The world report on violence and health”, *Lancet*, 2002. doi: [10.1016/S0140-6736\(02\)11133-0](https://doi.org/10.1016/S0140-6736(02)11133-0).
- [51] C. H. Demarty *et al.*, “The MediaEval 2013 affect task: Violent Scenes Detection”, in *CEUR Workshop Proceedings*, 2013.
- [52] T. Giannakopoulos, A. Makris, D. Kosmopoulos, S. Perantonis, and S. Theodoridis, “Audio-visual fusion for detecting violent scenes in videos”, in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2010. doi: [10.1007/978-3-642-12842-4_13](https://doi.org/10.1007/978-3-642-12842-4_13).

- [53] A. Ali and N. Senan, “Violence video classification performance using deep neural networks”, *Advances in Intelligent Systems and Computing*, vol. 700, pp. 225–233, 2018. doi: [10.1007/978-3-319-72550-5_22](https://doi.org/10.1007/978-3-319-72550-5_22).
- [54] T. W. Chua, K. Leman, and F. Gao, “Hierarchical audio-visual surveillance for passenger elevators”, in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2014. doi: [10.1007/978-3-319-04117-9_5](https://doi.org/10.1007/978-3-319-04117-9_5).
- [55] M. Bautista-Duran *et al.*, “Acoustic detection of violence in real and fictional environments”, in *ICPRAM 2017 - Proceedings of the 6th International Conference on Pattern Recognition Applications and Methods*, 2017. doi: [10.5220/0006195004560462](https://doi.org/10.5220/0006195004560462).
- [56] WHO. Department of Reproductive Health Research. London School of Hygiene and Tropical Medicine. South African Medical Research Council., *WHO | Global and regional estimates of violence against women*. 2013.
- [57] European Union Agency for Fundamental Rights, *Violence against women : An EU-wide survey*. 2014. doi: [10.2811/62230](https://doi.org/10.2811/62230).
- [58] K. Beyer, A. B. Wallis, and L. K. Hamberger, “Neighborhood Environment and Intimate Partner Violence: A Systematic Review”, *Trauma, Violence, and Abuse*, 2015. doi: [10.1177/1524838013515758](https://doi.org/10.1177/1524838013515758).
- [59] V. Mapell, *UPC-TALP database of isolated meeting-room acoustic events*, 2012. [Online]. Available: <http://catalog.elra.info/en-us/repository/browse/ELRA-S0268/>.
- [60] P. Foggia, N. Petkov, A. Saggese, N. Strisciuglio, and M. Vento, “Reliable detection of audio events in highly noisy environments”, *Pattern Recognition Letters*, 2015. doi: [10.1016/j.patrec.2015.06.026](https://doi.org/10.1016/j.patrec.2015.06.026).
- [61] E. Fagerlund and A. Hiltunen, *TUT Rare sound events*, 2017.
- [62] D. Stowell and E. Benetos, *IEEE AASP Challenge: Detection and Classification of Acoustic Scenes and Events*, 2013.
- [63] E. Cakir and T. Heittola, *TUT-SED Synthetic*, 2016.
- [64] C. H. Demarty, C. Penet, M. Soleymani, and G. Gravier, “VSD, a public dataset for the detection of violent scenes in movies: design, annotation, analysis and evaluation”, *Multimedia Tools and Applications*, 2015. doi: [10.1007/s11042-014-1984-4](https://doi.org/10.1007/s11042-014-1984-4).
- [65] J. F. Gemmeke *et al.*, “Audio Set: An ontology and human-labeled dataset for audio events”, in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 2017. doi: [10.1109/ICASSP.2017.7952261](https://doi.org/10.1109/ICASSP.2017.7952261).
- [66] E. Fonseca *et al.*, “Freesound datasets: A platform for the creation of open audio datasets”, in *Proceedings of the 18th International Society for Music Information Retrieval Conference, ISMIR 2017*, 2017.

- [67] Sound Understanding group, *AudioSet*, 2017. [Online]. Available: <https://research.google.com/audioset/index.html>.
- [68] M. A. Hearst, “Automatic acquisition of hyponyms from large text corpora”, 1992. doi: [10.3115/992133.992154](https://doi.org/10.3115/992133.992154).
- [69] A. Singhal, *Introducing the Knowledge Graph: things, not strings*, 2012.
- [70] GoogleResearch, “TensorFlow: Large-scale machine learning on heterogeneous systems”, *Google Research*, 2015. arXiv: [arXiv:1603.04467v2](https://arxiv.org/abs/1603.04467v2).
- [71] S. Hershey *et al.*, “CNN architectures for large-scale audio classification”, in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 2017. doi: [10.1109/ICASSP.2017.7952132](https://doi.org/10.1109/ICASSP.2017.7952132). arXiv: [1609.09430](https://arxiv.org/abs/1609.09430).
- [72] S. Abu-El-Haija *et al.*, “YouTube-8M: A Large-Scale Video Classification Benchmark”, *CoRR*, vol. abs/1609.0, 2016. [Online]. Available: <http://arxiv.org/abs/1609.08675>.
- [73] L. Deng, *A tutorial survey of architectures, algorithms, and applications for deep learning*, 2014. doi: [10.1017/atsip.2013.9](https://doi.org/10.1017/atsip.2013.9).
- [74] W. Di, A. Bhardwaj, and W. Jianing, *Deep Learning Essentials*. 2018.
- [75] A. Karpathy, *CS231n Convolutional Neural Networks for Visual Recognition*, 2016. [Online]. Available: <http://cs231n.github.io/convolutional-networks/>.
- [76] S. Saha, *No A Comprehensive Guide to Convolutional Neural Networks — the ELI5 way*, 2018.
- [77] A. Anwar, *Difference between Local Response Normalization and Batch Normalization*, 2019. [Online]. Available: <https://towardsdatascience.com/difference-between-local-response-normalization-and-batch-normalization-272308c034ac>.
- [78] Keras, *Convolutional Layers*.
- [79] H. Mahmood, *The Softmax Function, Simplified*, 2018.
- [80] T. Hinz, P. Barros, and S. Wermter, “The Effects of Regularization on Learning Facial Expressions with Convolutional Neural Networks”, 2016, pp. 80–87. doi: [10.1007/978-3-319-44781-0_10](https://doi.org/10.1007/978-3-319-44781-0_10).
- [81] K. J. Piczak, “Environmental sound classification with convolutional neural networks”, in *IEEE International Workshop on Machine Learning for Signal Processing, MLSP*, 2015. doi: [10.1109/MLSP.2015.7324337](https://doi.org/10.1109/MLSP.2015.7324337).
- [82] A. Kumar and B. Raj, “Deep CNN Framework for Audio Event Recognition using Weakly Labeled Web Data”, 2017.
- [83] ImageNet, *Results for ILSVRC2014*, 2014.

- [84] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition”, in *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*, 2015. arXiv: [1409.1556](https://arxiv.org/abs/1409.1556).
- [85] D. Ellis, S. Hershey, A. Jansen, and M. Plakal, *VGGish code*, 2017. [Online]. Available: <https://github.com/tensorflow/models/tree/master/research/audioset>.
- [86] H. Abdi and L. J. Williams, *Principal component analysis*, 2010. doi: [10.1002/wics.101](https://doi.org/10.1002/wics.101).
- [87] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, “SMOTE: Synthetic minority over-sampling technique”, *Journal of Artificial Intelligence Research*, 2002. doi: [10.1613/jair.953](https://doi.org/10.1613/jair.953). arXiv: [1106.1813](https://arxiv.org/abs/1106.1813).
- [88] J. Browniee, *SMOTE Oversampling for Imbalanced Classification with Python*, 2020.
- [89] R. Blagus and L. Lusa, “SMOTE for high-dimensional class-imbalanced data”, *BMC Bioinformatics*, 2013. doi: [10.1186/1471-2105-14-106](https://doi.org/10.1186/1471-2105-14-106).
- [90] Y. Xie, Y. Liu, and Q. Fu, “Imbalanced Data Sets Classification Based on SVM for Sand-Dust Storm Warning”, *Discrete Dynamics in Nature and Society*, 2015. doi: [10.1155/2015/562724](https://doi.org/10.1155/2015/562724).
- [91] S. Bhattacharyya, *Support Vector Machine: Kernel Trick; Mercer’s Theorem*, 2018. [Online]. Available: <https://towardsdatascience.com/understanding-support-vector-machine-part-2-kernel-trick-mercers-theorem-e1e6848c6c4d>.
- [92] G. Drakos, *Support Vector Machine vs Logistic Regression*, 2018.
- [93] C. Olah, *Understanding LSTM Networks*, 2015. [Online]. Available: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>.
- [94] SuperDataScience Team, *Recurrent Neural Networks (RNN) - The Vanishing Gradient Problem*, 2018. [Online]. Available: <https://www.superdatascience.com/blogs/recurrent-neural-networks-rnn-the-vanishing-gradient-problem/>.
- [95] M. Nguyen, *Illustrated Guide to LSTM’s and GRU’s: A step by step explanation*, 2018. [Online]. Available: <https://towardsdatascience.com/illustrated-guide-to-lstms-and-gru-s-a-step-by-step-explanation-44e9eb85bf21>.
- [96] M. Levoy, K. Dektar, and A. Adams, *Spatial Convolution*, 2012. [Online]. Available: <https://graphics.stanford.edu/courses/cs178/applets/convolution.html>.
- [97] M. Rampurawala, *Classification with TensorFlow and Dense Neural Networks*, 2019.

- [98] S. Kaski and J. Peltonen, “Dimensionality reduction for data visualization”, *IEEE Signal Processing Magazine*, 2011. doi: [10.1109/MSP.2010.940003](https://doi.org/10.1109/MSP.2010.940003).
- [99] J. Amat Rodrigo, *Análisis de Componentes Principales (Principal Component Analysis, PCA) y t-SNE*, 2017.
- [100] G. Hinton and S. Roweis, “Stochastic neighbor embedding”, in *Advances in Neural Information Processing Systems*, 2003.
- [101] L. Van Der Maaten and G. Hinton, “Visualizing data using t-SNE”, *Journal of Machine Learning Research*, 2008.
- [102] M. Wattenberg, F. Viegas, and I. Johnson, “How to Use t-SNE Effectively”, *Distill*, 2016. doi: [10.23915/distill.00002](https://doi.org/10.23915/distill.00002). [Online]. Available: <http://distill.pub/2016/misread-tsne>.
- [103] A. Kaw, *Holistic Numerical Methods*. [Online]. Available: <https://nm.mathforcollege.com/blog/>.
- [104] G. L. Prajapati and A. Patle, “On performing classification using SVM with radial basis and polynomial kernel functions”, in *Proceedings - 3rd International Conference on Emerging Trends in Engineering and Technology, ICETET 2010*, 2010. doi: [10.1109/ICETET.2010.134](https://doi.org/10.1109/ICETET.2010.134).
- [105] J. Browniee, *Difference Between a Batch and an Epoch in a Neural Network*, 2018. [Online]. Available: <https://machinelearningmastery.com/difference-between-a-batch-and-an-epoch/>.
- [106] Y. Wang, M. Huang, L. Zhao, and X. Zhu, “Attention-based LSTM for aspect-level sentiment classification”, in *EMNLP 2016 - Conference on Empirical Methods in Natural Language Processing, Proceedings*, 2016. doi: [10.18653/v1-d16-1058](https://doi.org/10.18653/v1-d16-1058).
- [107] J. Read, B. Pfahringer, G. Holmes, and E. Frank, “Classifier chains for multi-label classification”, *Machine Learning*, 2011. doi: [10.1007/s10994-011-5256-5](https://doi.org/10.1007/s10994-011-5256-5).
- [108] Indeed, *Salarios para empleos de Investigador/a en España*, 2020.
- [109] ——, *Salarios para empleos de Jefe de proyecto en España*, 2020.
- [110] Scikit-learn, *Metrics and scoring: quantifying the quality of predictions*. [Online]. Available: https://scikit-learn.org/stable/modules/model_selection/_devaluation.html.
- [111] A. Mishra, *Metrics to Evaluate your Machine Learning Algorithm*, 2018. [Online]. Available: <https://towardsdatascience.com/metrics-to-evaluate-your-machine-learning-algorithm-f10ba6e38234>.
- [112] F. Krüger, “Activity, Context, and Plan Recognition with Computational Causal Behaviour Models”, *ResearchGate*, 2018.

- [113] J. Browniee, *A Gentle Introduction to k-fold Cross-Validation*, 2018. [Online]. Available: <https://machinelearningmastery.com/k-fold-cross-validation/>.
- [114] S. M, *Why and how to Cross Validate a Model?*, 2018. [Online]. Available: <https://towardsdatascience.com/why-and-how-to-cross-validate-a-model-d6424b45261f>.
- [115] Scikit-learn, *Cross-validation: evaluating estimator performance*.
- [116] ML Glossary, *Loss Functions*, 2017. [Online]. Available: https://ml-cheatsheet.readthedocs.io/en/latest/loss%7B%5C_%7Dfunctions.html.
- [117] N. Laskaris, *How to apply machine learning and deep learning methods to audio analysis*, 2019.
- [118] Y. Lei, *Intelligent fault diagnosis and remaining useful life prediction of rotating machinery*. 2016, p. 37.
- [119] N. Kehtarnavaz, *Digital Signal Processing System Design*. 2008, pp. 175–196.
- [120] R. X. Gao and R. Yan, “Non-stationary signal processing for bearing health monitoring”, *International Journal of Manufacturing Research*, 2006. doi: [10.1504/IJMR.2006.010701](https://doi.org/10.1504/IJMR.2006.010701).
- [121] D. Gartzman, *Getting to Know the Mel Spectrogram*, 2019. [Online]. Available: <https://towardsdatascience.com/getting-to-know-the-mel-spectrogram-31bca3e2d9d0>.
- [122] L. L. Beranek and W. A. Rosenblith, “Acoustic Measurements”, *Physics Today*, 1950. doi: [10.1063/1.3066788](https://doi.org/10.1063/1.3066788).

A. APPENDIX METRICS

A fundamental part of a machine learning project consists of evaluating the performance. There are plenty of metrics to carry out this evaluation and the results will look in one way or another depending on the method utilized. The following two are the most used in this project

Classification Accuracy

This is a technique commonly used and it is usually referred to as accuracy. It can be defined as the relationship between the amount of right predictions and the total number on input instances [110].

$$acc = \frac{\text{Number of incorrect predictions}}{\text{Number of total input instances}}$$

This metric is appropriate when dealing with a balanced dataset, i.e., the same of number of samples per class. If the problem being addressed deals with unbalanced data, then the accuracy value could be a higher value due to predict all the instances belong to the major class. For example, if 90% of the data are part of the same class A and the model predictions are for this class, then the accuracy value will be 90%, which apparently is a satisfying output, even though we are misclassifying all the samples from class B [111] and the system is not learning.

Confusion matrix

As it own name indicates, the output of this type of metric consists on a matrix that shows a complete evaluation of the model. By definition, an entry i, j of the matrix denotes the amount of observations that belong to group i but are predicted as group j [110]. For example, considering a binary classification problem in which there are two classes, YES

$n = 165$	Predicted: NO	Predicted: YES
Label: NO	50	10
Label: YES	5	100

Fig. A.1. Example of confusion matrix

and NO, for a test set composed by 165 samples, the matrix included in figure A.1 is obtained.

There are four groups that can be extracted from this matrix: True positives, the samples that are predicted as YES and that is in fact their true label, True Negatives, those cases that were predicted as NO and they are originally labelled as NO, False Positives, in which the predicted label is YES but they are actually negative, and False Negatives, those in which the predicted label is NO when their original label is YES.

This metric an the one explained before, accuracy, can be related by taking the diagonal of the matrix and computing the next operation:

$$acc = \frac{TruePositives + FalseNegatives}{Total\ number\ of\ samples} = \frac{100 + 50}{165} = 0.91$$

When the classification task consists on more than two classes, a multiclass problem, a similar definition of the confusion matrix can be extended from the binary problem. Considering a certain observation C_k , the True positive part of the matrix is placed in the exact point where the column and the row of this certain observation are crossed, i.e, when the predicted label is equal to the true label. The False positives samples are placed along the column C_k for all the rows $C_0, \dots, C_{k-1}, C_k + 1, \dots, C_n$ which refers to all the samples that have been misclassified as the class C_k . The False negatives are, however, all the samples that originally are originally labelled as C_k but have been wrongly categorized as $C_0, \dots, C_{k-1}, C_k + 1, \dots, C_n$ classes. Finally, the True negative samples are distributed across all the other positions in the matrix. In figure A.2, an example for this explanation is shown.

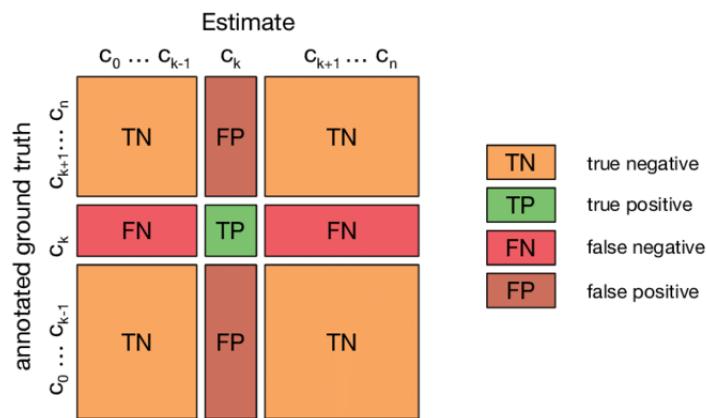


Fig. A.2. Confusion matrix for a multiclass classification [112]

B. APPENDIX K-FOLD CROSS VALIDATION

The cross-validation technique is a resampling practice that is commonly used in order to evaluate machine learning algorithms when the dataset is not very large. It just depends on the parameter k which represents the amount of folds or groups the data is going to be split into. This is the reason of the k-fold prefix. When the method refers to a situation in which the parameter is already fixed, for example, if $k = 5$, it becomes a 5-fold cross-validation.

The purpose of this procedure is to check the performance of a machine learning model on unseen data. This is done to avoid evaluating with the data used for the training process. It is highly used nowadays and allows to obtain a less biased estimation than that obtain by using other kinds of techniques as just a simple train and test split [113]. The steps followed by the method are listed below [114].

1. The data is split in a random way into k folds. The value of this parameter can be chosen previously by running the algorithm for different options and picking the one with the best performance. However, a not very high value is usually chosen, keeping it in a range between 5 and 10.
2. The model is fitted by using the data in the $k - 1$ folds as the train set, and the other part in the k fold for the validation process. Then the results are collected for this configuration.
3. The process must be repeated until all folds have been used as the validation set. Then, the average and standard deviation of all the results can be computed as the metric of the model.

Sometimes, instead of dividing the whole set of data in k folds, it is first done a split into train and test sets. Then, the test set is excluded for a final measure and the train is split again with the k-Fold Cross Validation procedure. This is a good choice when data augmentation techniques that alter the a priori distribution of the data are employed. In this case, the test set retains the original a priori class distribution. The concept of the technique is shown in figure B.1.

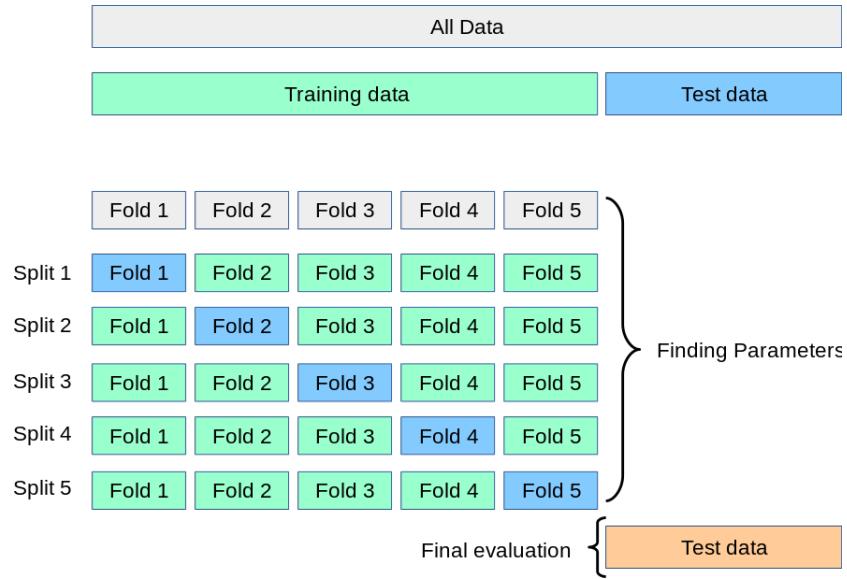


Fig. B.1. K-fold cross-validation scheme. The data is first split into train and test, and then the train is split again with this method [115]

C. APPENDIX CATEGORICAL CROSS-ENTROPY

In order to define this kind of loss function, it is necessary to consider the *softmax* definition. As it was already explained in 3.2.1, this function takes a vector and represents its elements in the range (0,1) as probabilities, so that all the resulting values must sum up to 1. In a multiclass classification, these are interpreted as class probabilities.

These output probabilities can be evaluated with the cross-entropy loss. This is a way of measuring the performance of a model in which the output are probabilities between 0 and 1. The loss increases if the probability differs largely from the actual value of the true label. A perfect system would have a loss value of 0 [116]. In figure 5.10, an example of the loss function when the true label is 1 is shown.

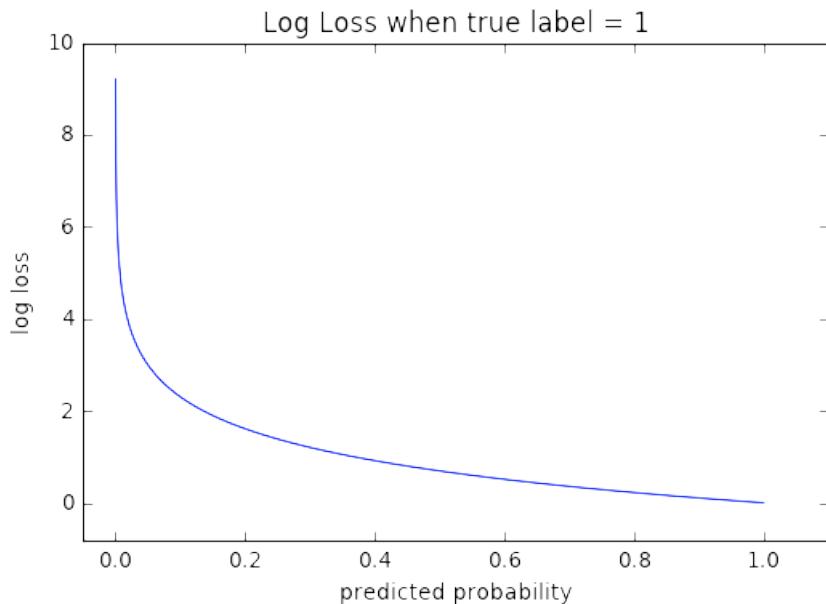


Fig. C.1. Cross-entropy loss function when the true label is equals to 1. As the probability of the predicted class approaches to 1, the loss function tends to zero. However, if it the probability is closer to 0.0, then the loss function increases heavily [116].

Its formula is defined for binary problems as follows:

$$CE = -(y \log(p) + (1 - y) \log(1 - p))$$

For a multiclass case, a unique loss is computed for each class label for each sample, and then the results are added as follows:

$$CE = - \sum_{c=1}^M y_{o,c} \log(p_{o,c})$$

Where M is the number of classes, y is a indication in binary format if the class c is the right prediction for this sample o and p is the estimation probability of c for o [116].

Then, the categorical cross-entropy is a combination of softmax and the cross-entropy loss.

D. APPENDIX SPECTROGRAM AND MEL SCALE

In a natural situation the way humans understand audio as we hear it, can be interpreted as a succession of sounds that evolve with time. Actually, what happens during a recording process is that the signal registered is nothing but the variation of the values that represent the changes in the acoustic pressure along time, what is usually called the amplitude of the signal. This is what we see when visualizing a raw audio file in its original form, i.e. in the time-domain [117], as in figure D.1.

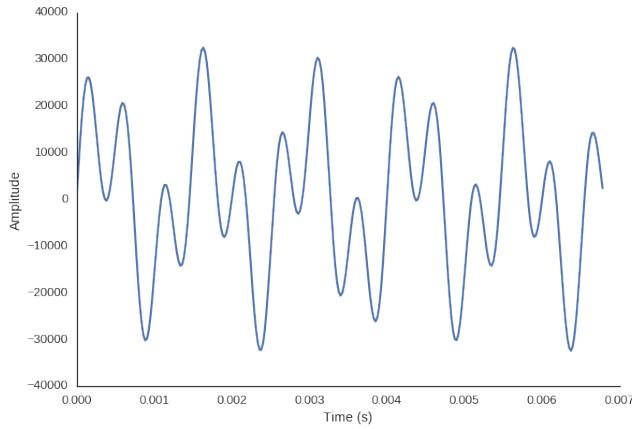


Fig. D.1. Audio signal in time domain

However, in this form, audio data is not very useful for learning tasks. With the objective of representing audio data in another format so information can be extracted in a more practical way, it is normally converted to the frequency domain. The FT is the technique in charge of performing this conversion. Since in signal processing we work with discrete data, the algorithm is called Discrete Fourier Transform (DFT), which is essentially the same as FT but for discrete data and frequencies. The most common algorithm to compute this is the one called Fast Fourier Transform (FFT), which performs the DFT in a more efficient way [118]. The result of this method shows the values of the amplitude against the frequency axis. However, when the input is not stationary but it still can be regarded as such when observed within small portions, there is also another method to perform this operation and it is called the Short-Time Fourier Transform (STFT), which is the one used in this work.

The STFT can be defined as a sequence of FTs which are computed along a windowed signal. So, apart from computing the transformation, it also localizes in time this frequency information, given as a result a variation of the frequency along the time axis. Below, the formula of this operation is included:

$$X_{STFT} = \sum_{k=0}^{L-1} x[k]g[k-m]e^{-j2\pi nk/L}$$

$x[k]$ is the time signal of length L and $g[k]$ is the window that goes through all the signal. In figure D.2, a visualization of the whole process is shown [119].

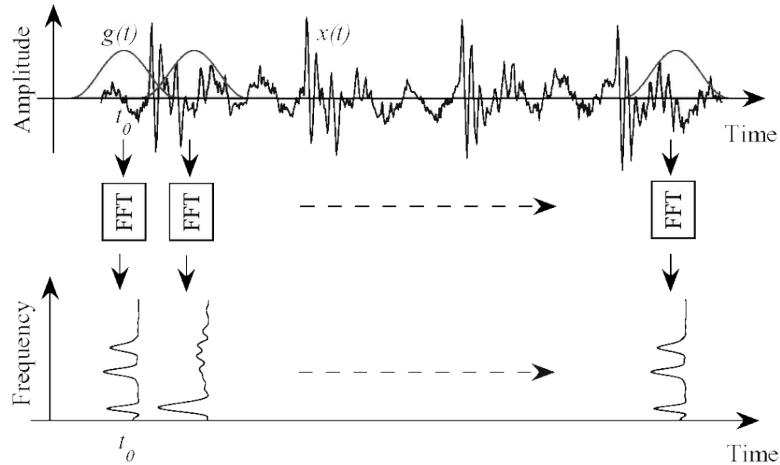


Fig. D.2. STFT [120]

What is obtained from this computation after transversing the whole signal, is a representation that depends on three magnitudes: time, frequency and the amplitude. This is known as the spectrogram of the signal. With this visualization we can see how the energy of the signal varies with time and also how it is placed in the different frequencies. In figure D.3 it is shown how a typical representation of a spectrogram looks like.

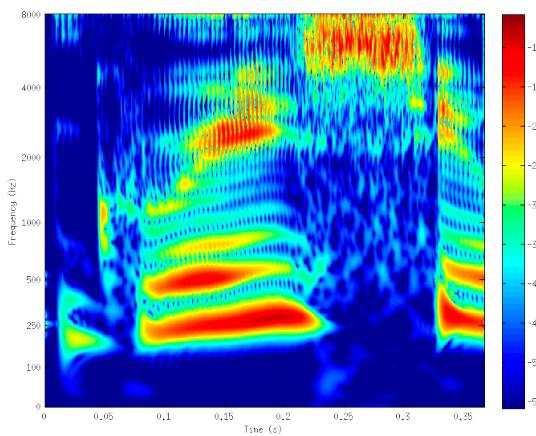


Fig. D.3. Spectrogram of an audio signal. The x-axis corresponds to time, the y-axis to frequency and the colour of the plot to the amplitudes by following the colour bar on the right.

When dealing with audio learning tasks we usually employ a logarithmic transformation of the frequency domain such as the Mel scale. It has been successfully employed

in plenty of applications because its similarities with how humans process sounds [121]. This is done by considering the fact that we are more capable of differentiating small variations of pitch at low frequencies than at high ones. The following formula is the way to convert from frequency to mel scale [122]:

$$M(f) = 1127 \ln\left(1 + \frac{f}{700}\right)$$

This same idea of applying the human way of understanding sounds has also been used to generate the Mel-frequency Cepstrum Coefficients (MFCC), which have been widely used in speech related problems. The first thing to understand in this area is that the sounds we produce by the lungs that make the vocal folds vibrate are subsequently filtered by the vocal tract. If a similar representation could be modelled in an artificial system, a more accurate depiction of the sound uttered could be achieved. Actually, the envelope of the short time power spectrum of the signal can be considered as the shape of this filter and it can be reliably represented by the MFCC, so useful information can be deduced from them. This type of features have been extensively used for several fields such as automatic speech generation or speaker recognition but also in audio related ones [20].