



Deception Detection Using Facial Action Units

Facoltà di Ingengeria dell'informazione, Informatica e Statistica
Corso di Laurea Magistrale in Informatica

Candidate

Emanuele Orfanelli
ID number 1383726

Thesis Advisor

Prof. Luigi Cinque

Co-Advisor

Dr. Danilo Avola

Academic Year 2018/2019

Deception Detection Using Facial Action Units
Master's thesis. Sapienza – University of Rome

© 2018 Emanuele Orfanelli. All rights reserved

This thesis has been typeset by L^AT_EX and the Sapthesis class.

Author's email: emanueleorfanelli@gmail.com

*Dedicated to
my Family and Friends*

Acknowledgments

I want to thank my family for the support they gave me over the years.

A special thanks to my friends who have always been there for me, especially Enrico de Francisci, Francesco Andreozzi, Cabio Ficcarelli.

I'd like to thank Prof. Luigi Cinque for trusting me with this work and for the availability he showed me. I'd also like to thank Dott. Danilo Avola for helping me through out this thesis.

Contents

1	Introduction	1
1.1	Overview of the work	1
1.2	People Lie Detection	3
1.3	State of the Art	5
1.3.1	Speech	5
1.3.2	Eyes	7
1.3.3	Neuroscience	8
1.3.4	Thermal Imaging	10
1.3.5	Multimodal	11
1.3.6	Facial Expression and Micro-Expressions	13
1.3.7	Action Units	15
1.4	Contributions and Outline	18
1.4.1	Contribution	18
1.4.2	Outline	18
2	Machine Learning	19
2.1	Introduction	19
2.1.1	Supervised Learning	20
2.1.2	Unsupervised Learning	23
2.1.3	Classification	24
2.1.4	Regression Analysis	26
2.1.5	Linear Regression	28
2.2	Logistic Regression	30
2.3	Random Forest	33
2.4	SVM	40
2.4.1	Hyperplane	40
2.4.2	Maximal Margin Classifier	41
2.4.3	Support Vectors	42
2.4.4	Finding the Maximal Margin Classifier	42
2.4.5	Support Vector Classifier	43
2.4.6	Kernels	46
3	Architecture	49
3.1	General Architecture	49
3.2	OpenFace	50
3.3	Landmark Detection	51

3.3.1	CEN	52
3.3.2	Point Distribution Model	53
3.4	Action Unit Detection	56
3.4.1	Training Datasets	56
3.4.2	Feature Extraction	58
3.5	Deception Detection	62
3.5.1	Data pre-processing	62
3.5.2	SVM	63
4	Experiments	65
4.1	Building Blocks	65
4.2	Real Life Trial Database	66
4.3	Data Analysis	69
4.3.1	Data Comparison	69
4.3.2	Correlation	71
4.3.3	GLM	72
4.3.4	Random Forest	73
4.3.5	SVM	74
4.4	Results	75
5	Conclusions	77
5.1	Final Considerations	77
5.2	Future Work	77

Chapter 1

Introduction

In this chapter we give an overview of the work (Par. 1.1). We start examining how people perform at detecting lies in order to have a comparison with other methods (Par. 1.2).

We then present a taxonomy of the current state of the art (Par. 1.3) concerning lie detection with particular emphasis on computer vision.

The last section is about the structure of the rest of the thesis, and our contribution to the field of deception detection (Par. 1.4).

1.1 Overview of the work

Deception detection has always been a very interesting topic since it has numerous social and psychological implications in many different fields. In the past 40 years there has been a steady increase in researchers interested in studying deceptive behavior, first from a sociological and psychological point of view, and in more recent years from a technological standpoint, aided by the advance of machine learning techniques, especially in the field of Computer Vision, an interdisciplinary field that aims to simulate the human visual system to automatically comprehend images or videos.

Our work aims to recognize whether a person is lying or telling the truth by performing a frame by frame analysis on a database of videos, with the goal of extracting a set of the different facial movements performed by the person in the video. After finishing the extraction, the data is submitted to analysis by using statistical and machine learning techniques.

The dataset we are using was provided by Perez-Rosas et al. [71] and is composed by 121 videos taken from real life trials and uploaded on YouTube (Par. 4.2). Since those videos are taken from real life events, the illumination, pose, audio and video features are not homogeneous and often substandard. Substantial work was done to eliminate the parts that are not relevant, and as result some videos were trimmed or discarded.

After trimming and cleaning the videos, we split the data into two sets (train and test) based on the subject ID and extracted all the facial movements. Then we utilized machine learning and statistical approaches to understand if the person in the video was lying, based on the results of analyzing each frame.

The result of this work, especially when improved by having a better and larger dataset, can be useful in different situations, even though it is important to remember that privacy is a real concern, and with this kind of applications it should never be underestimated.

Some of the applications of this work can include airport security, work interviews, and possibly many kinds of social interactions. It could be eventually used by the public to review a speech of a political candidate or by the police force as an aid in the interrogation process. Another use can be in trials where people lives are at stake and discerning a true or false testimony might be vital. We really hope this work proves to be socially useful.

1.2 People Lie Detection

It's very rare for people to be able to consistently discern between lies and truths, even though we hear lies regularly in the course of our lives. In fact most untrained people perform like chance (50%) when tasked with detecting lies [73], and considering that the ordinary person lies at least twice a day on average [58], and that in this digital age the amount of lies told daily are increasing [41] due to on-line communication making us feel more protected and confident in our deceptions, the problem of detecting lies is an important one and it deserves a good deal of research effort.

People's difficulty in detecting lies stem from the lack of objectivity. We are biased by so many factors and skilled deceivers can take advantage of that.

An important reason is that people generally think it is easy to spot liars, underestimating the effort it takes by having too much confidence on their own judgment and capabilities, and by believing in gut feelings instead of empirical clues [75].

In a study by Baker et al. [9] 116 participants were asked to judge 20 videos of people pleading for the safe return of missing family members, half of which were lies where the pleader was the one responsible for the disappearance of the victim. The participants provided confidence ratings, the cues they utilized to make their judgment, and their emotional response to each video.

Weirdly, emotionally intelligent people perform worse at deception detection. This is due to their greater sympathetic feelings towards others (enhanced gullibility) that can often cloud their judgment.

In [20] De Paulo et al. analyze the accuracy of deception judgments from a collection of 206 documents and 24.483 judges. They found that people can differentiate lies and truths with an accuracy of 54%, with a lie detection accuracy of 47% and a truth detection accuracy of 61%. Interestingly, their findings reveal that is easier for people to discriminate between lies and truths from the audio cues rather than the visual ones.

In another study [43] 192 students obtained an accuracy of 55.2% on lie detection, with 61% accuracy for guilty suspects and 49% for innocent ones.

Police officers and other trained officials seems to perform better than the average person at detecting lies [97], obtaining around 70% accuracy for distinguishing both lies and truths correctly, supporting the hypothesis that training and practice can make us better at spotting liars.

<i>Lie Catchers</i>	<i>Number of Lie Catchers</i>	<i>Overall Accuracy</i>	<i>Liar Spotting Accuracy</i>	<i>Truth- teller Spotting Accuracy</i>
Diverse (206 documents synthesized)	N/A	54%	47%	61%
Undergraduate students	192	55.2%	61.5%	49%
Police officers, CIA and FBI agents, lawyers, college students, therapists, judges, etc.	N/A	50%	N/A	N/A
Secret Service agents	N/A	70%	N/A	N/A
Police officer	37	72%	73%	70%
College students	20	50%	N/A	N/A

Figure 1.1. Performance on deceit detection by human observers [90].

1.3 State of the Art

How are lies detected? At the moment there are a lot of different instruments and technologies to attempt to detect deceit, ranging from analyzing psychological features, to using the polygraph or thermal cameras and machine learning approaches.

In computer vision specifically, lie detection is a somewhat new field of research, and as such it's done using many different approaches, employing not only RGB cameras but also physical sensors and thermal cameras, and using voice, gaze, action and body analysis tools, all of this often aided by machine learning techniques, and by combining many of these modalities together to achieve better results.

We now proceed to describe the state of the art, based on the latest researches done in the field, with specific focus on computer vision:

1.3.1 Speech

Speech is one of the many methods that can be used to recognize if a person is lying. In fact, the speech signal contains linguistic, expressive, organic and biological data [69].

One of the most used indicator of deception in various studies has been the response latency [96], since inventing a lie requires additional cognitive load as opposed to remembering the truth. The authors also notice that habitual lying makes it easier, and conversely often being spontaneous and telling the truth makes lying harder. Another indicator of lying is the speech rate, especially when it's different from the average rate of a specific person [18]. Other verbal cues like grammar usage and word frequency have been used and have achieved high accuracy in psychological researches [74].

Speech analysis can reveal changes that affect behavior, such as stress, emotion, deception etc. by analyzing the pitch and the stress level. When a stressful situation arises, the hormonal levels of the body change, and this causes an increase in blood pressure and heart rate. This in turn affects the muscle in the respiratory system, and so speech is affected [69].

In sound processing, the mel-frequency cepstrum (MFC, Fig. 1.2) is a representation of the short-term power spectrum of a sound. Mel-frequency cepstral coefficients (MFCCs) are coefficients that collectively make up the MFC [103].

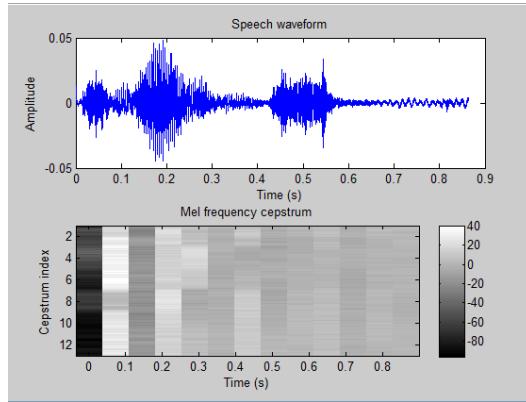


Figure 1.2. Topside is the spoken word, bottom is the MFCC derived from the word

In [64] the authors created a new database by making 40 candidates try to deceive them while telling truthful or deceptive statements for about one to two minutes each. From this experiment they extracted MFCC and pitch, so that they could process them through Matlab's Voice Box.

After acquiring the data, an SVM classifier was trained to classify new data, obtaining an accuracy of lie and truth detection from speech audio respectively of 88.23% and 84.52%.

In [63] Perez et al. utilize a dataset of 140 videos created by users who uploaded either truthful or deceptive statements. They collected the transcriptions from the videos both manually and via automatic speech recognition software. By using a bag-of-words representation of the transcripts they build a vocabulary, and to classify the data they employed an SVM and Naive Bayes classifier, achieving 73% accuracy.

1.3.2 Eyes

Using the eyes to detect lies is one of the most studied approaches, as the eyes hold a significant amount of information regarding our thinking and emotional state [36] (Fig 1.3). We will focus our description on the Computer Vision approaches rather than psychological ones.



Figure 1.3. The eyes could hold a lot of informations regarding what we are thinking [27].

An important advantage of using the eyes as a vehicle for deception detection is that it is possible to generate a non invasive approach while analyzing them, meaning that subjects do not need to willingly participate or even know if they are being examined or not (the moral matter should be considered in another setting), and there could be no need to have huge and expensive machinery, like for example the polygraph or a fMRI machine (par. 1.3.3).

Cognitive load, which is set to increase while lying, is one of the most significant factor for deception detection and can be analyzed through the eyes. Important factors are also blink rate, gaze aversion and pupil dilation.

In [38] the authors utilize high speed cameras to record and analyze blink count and blink duration of 50 subjects while asking 10 control questions, to see if there is a variation in them while the subject is being questioned. The authors analyzed the resulting images frame by frame and based on the facial landmarks around the eye they recognize AU45, the action unit (muscle movement) associated blinking. The results shows that both blink duration and count are increased while lying.

In another study, Leal and Vrij [53] gathered 26 people, 13 liars and 13 truth tellers, to ask them to lie or tell the truth in a target period, while having a baseline from two preceding periods. The eye blinks during the target and baseline periods and directly after the target period (target offset period) were recorded.

Compared to the baseline periods, lying subjects show a decrease in eye blinks during

the target period and an increase in eye blinks during the target offset period. This means that there is a decrease in eye blinks while lying and an increase just after lying. This pattern resulted very different for truth tellers, showing that there is a significant difference in eye blink behavior between truthfulness and deception.

Singh et al. in [86] show that while lying there is an increase in cognitive load and a significant decrease in eye blinks, directly followed by an increase in blink rate as soon as the cognitive demand ceases, after telling the lie. A threshold is set by the authors for this study, either at 26 blinks per minute or it is calculated ad personam using the average blink rate from a blink detection algorithm. Blink detection is done through MATLAB using the HAAR Cascade algorithm.

Lim et al. study eye gaze [56] to investigate the relation with lie detection. The result supports the theory that an increase in cognitive load leads to a decrement in the number of eye movements while lying.

Bhaskaran et al. measure deception by the deviation from normal behavior at critical points during an investigative interrogation [16]. For starters, a dynamic Bayesian model of the eye movement is trained during a normal conversation with each of the 40 subjects of the experiment.

The remainder of the conversation is then broken into pieces and each piece is tested against the normal behavior. The deviation from normality are observed during critical points in the interrogation and used to deduce the presence of deceit, obtaining an accuracy of 82.5%.

In [77] Proudfoot et al., using latent growth curve modeling, research how the pupil diameter changes over the course of an interaction with a deception detection system. The assumption is that anxiety changes the pupil diameter. The subjects are presented with crime-relevant target items (possibly incriminating) and non relevant items.

The results indicate that the trends in the changes are indicative of deception during the interaction: pupil diameter is initially bigger for guilty subjects relative to innocent ones. Also the pupil diameter of the participants decreased between subsequent sections of the test, specifically the magnitude of the decrease was greater for guilty subjects who did not see incriminating items.

1.3.3 Neuroscience

Electroencephalogram (EEG) and functional Magnetic Resonance Imaging (fMRI) have been employed for lie detection with decent results for a long time, but at the cost of invasiveness, since both methods require big machinery, a suitable environment, and subjects willing to participate by having electrodes and other monitoring tools attached to their head and body.

EEG is a monitoring method that records brain activities based on its potential.

In [85] Simbolon et al, the authors use Event Related Potentials (ERP) to measure brain response directly from thought or perception. Among the many types of signals that constitute the ERP signal, P300 (Fig. 1.4) is the most critical for lie detection,

as it appears as a response to meaningful rare stimuli (called odd ball stimuli). Eleven males of age between 20 and 27 took part in the study. The gathered data were then divided into training and test sets to produce different models. The highest accuracy of 70.83% was reached by a SVM classifier in detecting lying subjects.

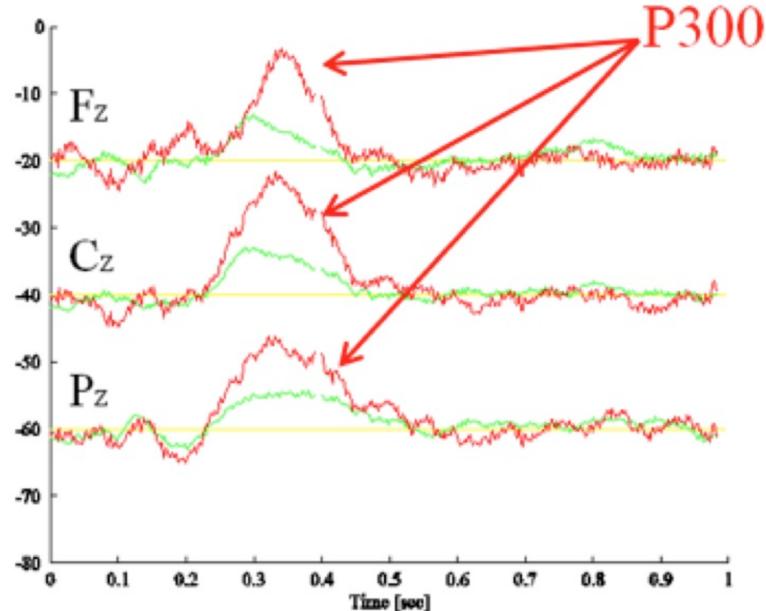


Figure 1.4. EEG of P300 waves image on channels Fz, Cz, and Pz

In the study performed in [52] twenty people of ages between 22 and 24 years old were subject to a card test using an EEG machine. EEG signals were collected using electrodes attached to the subject's head. The authors used the EEG signals to identify useful frequency bands and to measure the lying state based on spectral analysis, with the use of fuzzy reasoning, obtaining 89.5% detection accuracy.

Arasteh et al. [11] utilize an alternative approach to the polygraph, the P300 Guilty Knowledge Test (GTK). GTK is based on the amplitudes of the P300 ERP wave as an index for the subject's recognition of concealed information. The Guilty Knowledge Test works on the assumption that among many similar unfamiliar topics, the recognition of familiar ones will be followed by a different response.

62 subjects were part of this experiment and participated in a mock crime followed by the P300 Guilty Knowledge Test. The authors used empirical mode decomposition (EMD) to extract features from the EEG signal and modeled them through matlab. A genetic algorithm was then utilized for the feature selection and to handle the dimensionality increase of this approach. The classification accuracy of guilty and innocent subjects was 92.73%.

Another neuro-scientific approach to lie detection is functional Magnetic Resonance Imaging (fMRI). fMRI measures brain activity by detecting changes associated with blood flow (Fig. 1.5). This technique relies on the fact that cerebral blood flow and neuronal activation are coupled. When an area of the brain is in use, blood flow to that region also increases [100].

In most of the experiments with fMRI and lie detection, candidates are instructed to lie and tell the truth in specific situations, and the brain activity from these instances is compared to a baseline condition. The regions showing greater activation for lies than truth are supposed to be the most significant for deception detection. In a recent study it is shown that there is considerable agreement on the significant areas of the brain that regard lying [35].

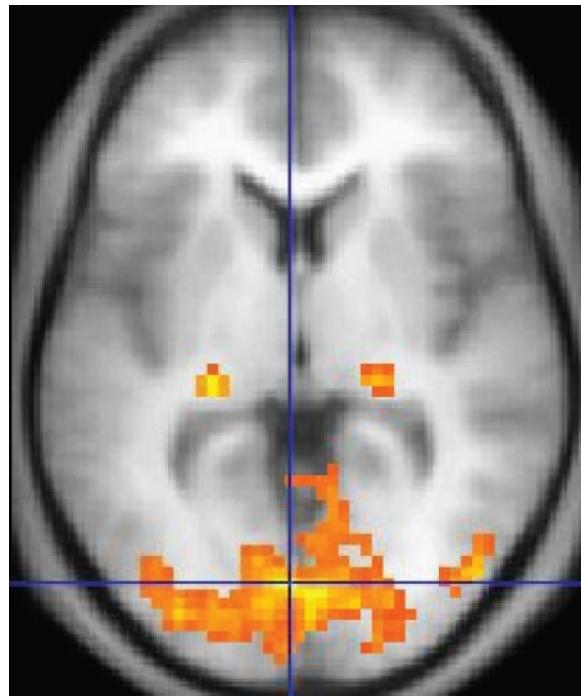


Figure 1.5. fMRI image with yellow areas showing increased activity compared with a control condition [100]

As much as fMRI is useful for lie detection, it presents some shortcomings [26] [3]: many fMRI studies are small, not replicated and done with just a few subjects; there are some contradicting results between some studies; most of the studies are not done in a context of high stake deception, but in a controlled environment where subjects are asked to lie about some topic or event, but often without a real interest in being deceitful. Another important point is that the fMRI approach requires collaboration and expensive equipment to be carried out, so this is a very limiting factor.

1.3.4 Thermal Imaging

In thermal imaging, thermal features are extracted from the face using a high definition thermal camera. The objective is to analyze what kind of differences occur when a subject responds truthfully or deceptively to particular stimuli.

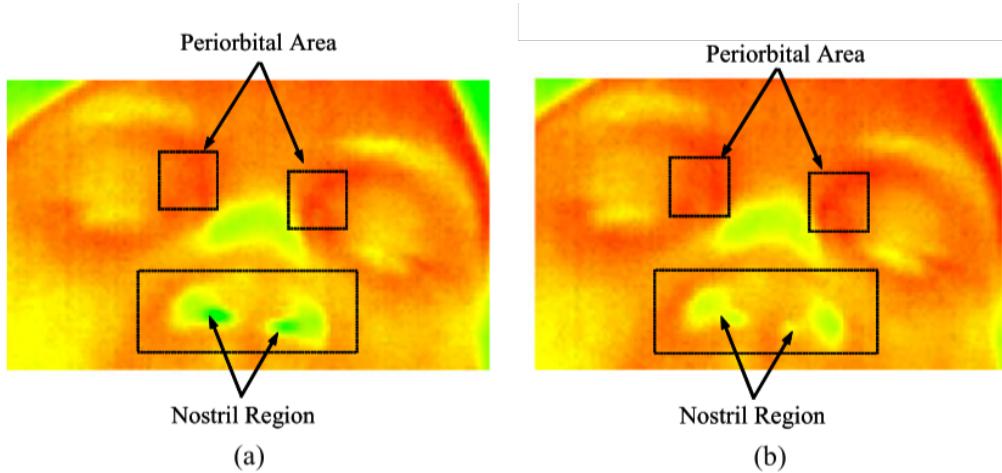


Figure 1.6. Examples of thermal images during (a) questioning and (b) answering [91]

According to a recent study [6] examining 30 subjects, the most relevant zones for deception detection, seen by utilizing thermal equipment, and located in the facial area are the forehead and periorbital regions.

In this study the subjects were registered with a thermal camera for one minute to extract the baseline features, and after that the interviewers asked a series of questions. A thermal map was created from the registrations using the Hue Saturation Value to differentiate between lies and truths.

In [79] the authors set up an experiment to collect 492 responses from 25 participants, using a deception scenario requiring the subjects to learn a story, provided by the authors, and practice their stories before an interview by giving them sample questions, so that the cognitive load would increase when receiving new and unseen questions during the interview process.

At the beginning of the interview four baseline questions were asked to register the initial thermal state of the subjects, and then a series of questions were proposed, with answers both present and missing from the previously provided story. After extracting the periorbital region's thermal variation, a k-nearest neighbor classifier was used, with an 87% accuracy in predicting lies or truths.

In [91] data are gathered non-intrusively from the nostril and periorbital regions using two dimensional far infrared cameras. The study lasted for two years and covered 18 test subjects involved in real crime cases. The temperature is extracted and converted in change in blood flow velocity, and a signature of the respiration pattern is determined in terms of the ratio of the measured maximum and minimum temperatures in the nostril area. The classification rate for this study is 88.5%.

1.3.5 Multimodal

After seeing all these different techniques to detect lies, it's only natural for researchers to try and fuse or aggregate some of them to see if it's possible to achieve better results [7].

In [5] Abouelenien et al. collect data from a dataset of 30 subjects to examine thermal and visual clues of deception. Their aim is to identify the regions that offer higher capability of detecting deceit.

The method employed uses the CERT (Computer Expression Recognition Toolbox) to detect facial expression and encodes them through Action Units (explained in greater detail in par. 1.3.7), a way to codify almost any kind of muscle movements on the face.

To extract thermal features they create a thermal map using gray-scale and Hue Saturation Value. They also calculated normalized blinking rates and the mean head orientation angle. In addition, over 60 physiological features are extracted and stored with the use of other sensors and RGB cameras.

The experimental results show that the non-contact feature fusion model outperforms traditional physiological measurements, and that the forehead region is one of the most promising areas to gather information for deception detection using thermal imaging.

In a following paper [4] Abouelenien et al. explore a multimodal deception detection approach comprised of physiological, linguistic, and thermal features on a new dataset of 149 recordings. They set out to determine which is the most discriminative region of the face to differentiate between truth and lies based on thermal imaging, and perform feature analysis using a decision tree model. The results show that the forehead could be a better indicator of deceit than the periorbital area. The physiological features did not contribute very much, while the linguistic feature played a critical role, where self-referencing and exaggeration words were the major indicators of deceit. The overall accuracy of the system is 70%.

Another example of multimodal detection is found in [108] where Wu et al. develop a framework to automatically detect deception in videos of trials. They utilize three modalities: vision, audio and text. For vision, they employ various classifiers trained on low level video features to predict human micro-expressions, and to successively predict deception. Interestingly, IDT (Improved Dense Trajectory) features, often used to recognize actions in videos, are good predictors of deceptive behavior.

The authors decided to fuse the score of the classifiers on IDT and micro-expressions to boost the performance. Regarding text, the transcript of the considered videos are analyzed, but the performance increase is very marginal.

For speech, they integrated the vision side with Mel Frequency Cepstral Coefficient (MFCC) features analysis from the audio, boosting the performances significantly, reaching an AUC of 0.877.

Noje et al. [68] set up a study with ten subjects to observe the potential of head movements in lie detection. They built an application to detect head movement and position by performing a frame to frame analysis on a video stream. A correlation was made between head movement/position and the identification of lies. The results of the study are not concluding as this data can't be utilized without being incorporated with other modalities such as voice, gaze, words, expressions et cetera.

In [71] Perez et al. utilize real life trial videos to identify deception, employing a model that extracts features from both linguistic and gesture modalities. The verbal features consist of unigrams and bigrams derived from the bag-of-words

representation of the video transcripts, while gesture features are taken from the annotations performed using the MUMIN coding scheme, where each feature indicates the presence of a gesture only if it's present in the majority of the video. Using Decision Trees and Random Forest classifiers they achieve between 60% and 75% accuracy.

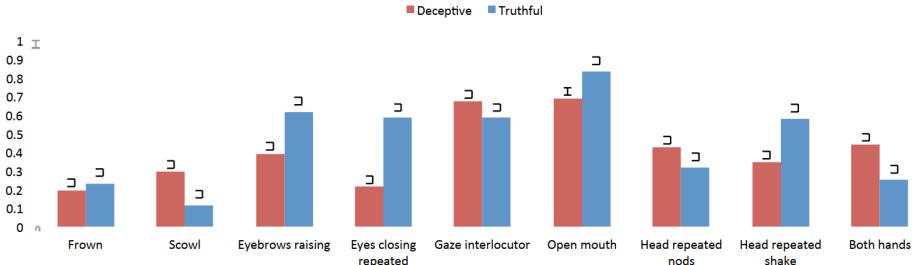


Figure 1.7. Distribution of non-verbal features for deceptive and truthful groups [71].

1.3.6 Facial Expression and Micro-Expressions

Facial Expressions are one of the main methods that we use to express our emotions, and are developed since the first months of our life to help us communicate our feelings to others. But what happens when we want to hide our emotions instead?

Facial micro-expressions are very fast (from 1/2 to 1/25 of a second) and involuntary expressions that come up on the human face when we are trying to suppress or hide an emotion, and are very difficult to control using just one's willpower [32].



Figure 1.8. Six basic Facial Expressions in adults and children [76].

Studying and classifying micro-expression is very valuable and has many applications, especially in psychology and forensic sciences, but it is a hard feat as the duration is

very short and the intensity is generally very low.

Micro-expressions have been studied since 1966 to recognize and distinguish real or fake emotions, initially by Haggard and Isaacs [39], and three years later by Ekman and Frisen [31].

Substantial work on Micro-Expressions has been done in recent years by Pfister, Li et al. In [72] they collaborated with psychologists to design an induced emotion suppression experiment. The data was collected with a high speed camera to be able to register the micro expressions of the subjects. A Temporal interpolation model was used to counter the shortness of the video length, while multiple kernel learning, random forest and SVM were used to classify micro expressions reaching an accuracy of 71.4%.

The lack of a big and well formed database was one of the biggest hindrance to their research. To solve this problem in [54] they reveal a new dataset, the Spontaneous Micro-expression Database (SMIC), which includes 164 micro-expression video clips taken from a group of 20 participants.

They used two high speed cameras to record the face of the subjects while they were watching a selection of videos that induced strong emotional responses, and they had to try and suppress those emotions. After the video the subjects had to answer about the emotions they felt while watching it. The data were then segmented and labeled by two annotators.

A study of spontaneous micro expression spotting and recognition methods was done in [55]. A new training-free method, based on feature difference contrast for recognizing micro-expressions was presented. The features are extracted from the video using Local Binary Pattern (LBP) and Histogram of Optical Flow.

To recognize the Micro expressions, the authors performed face alignment and temporal interpolation, then trained an SVM classifier.

This micro-expression framework was tested on the SMIC and CASMEII database with very good results. After combining micro-expressions spotting and recognition they released a new micro-expression analysis system (MESR) that is able to recognize micro-expressions from spontaneous video data.

Owayjan et al. [70] designed a lie detection system using micro-expressions. At first an embedded video system is used to record the subject interview. The video stream is converted into frames, and each frame is processed in four stages: converting the images, filtering out useless features, applying geometric templates and finally extracting the measurements to detect the micro-expressions.

Results show that up to eight facial expressions can be recognized, and that lies can be discerned with high precision.

In [49] Kawulok et al. explore how to exploit fast smile intensity detectors to extract temporal features using a SVM classifier. Using exclusively a face detector, without localizing or tracking facial landmarks, they analyze the smile intensity time series. They then employ an SVM classifier to improve training from weakly labeled datasets. Then, to train the smile detectors, they use uniform local binary pattern features. This allows to detect, in real time, between spontaneous or posed expressions.

Su et al. [90] aim to test the validity of facial clues to deception detection in

high-stakes situations using computer vision approaches. By using invariant 2D features from nine separate regions of the face they perform facial analysis on eye blink, eyebrow motion, wrinkle occurrence and mouth motion, integrated with a facial behavior pattern vector. Training a Random Forest to classify the patterns into deceptive or truthful, they achieved a 76.92% accuracy.

1.3.7 Action Units

Action Units (AUs) are defined as a contraction or relaxation of one or more muscles. They have been used in the Facial Action Coding System (FACS), a system developed initially by Hjorntsjö [45] and then improved by Freisen and Ekman [30], that categorizes all the facial movements by their appearance on the face.

Using FACS it's possible to code nearly any facial expression by deconstructing them into Action Units. For example the Duchenne smile (felt smile) is a combination of AU6 (orbicularis oculi, pars lateralis) and AU12 (zygomatic major) as we can see from fig. 1.9.

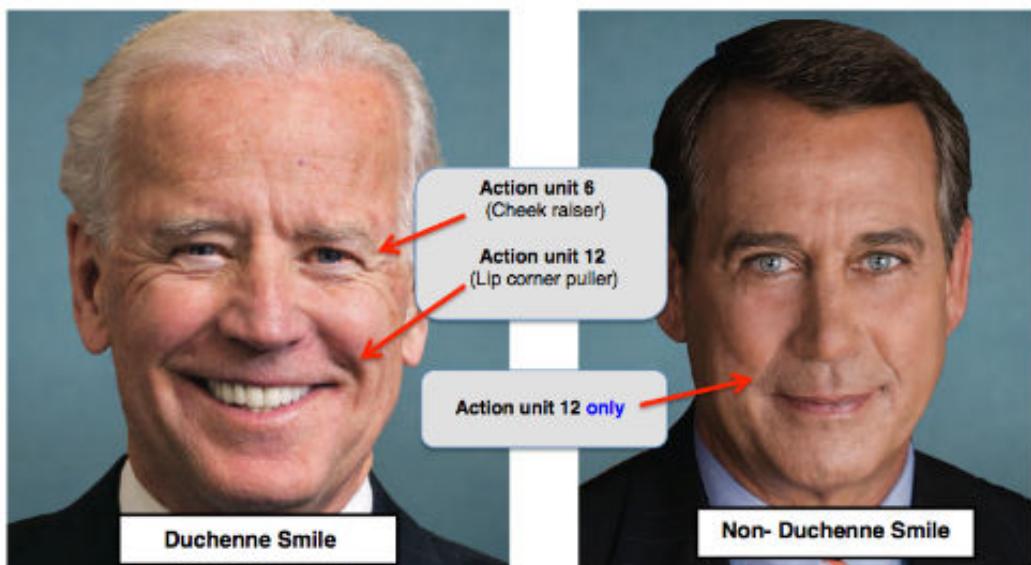


Figure 1.9. Duchenne smile (real) vs non Duchenne (fake).

In the FACS there are 44 AUs categorized with five level of intensity (A to E). Most of the AUs have an onset, peak and offset phase, as shown in fig. 1.10.

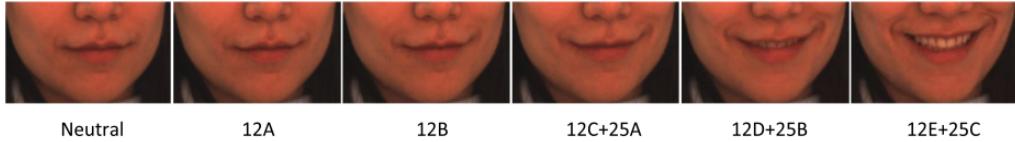


Figure 1.10. Facial images with AU intensity variations (AU12: lip corner puller, AU25: lips part) [60].

By using AUs is possible to recognize emotions based on the combination of AUs displayed (Fig. 1.11).

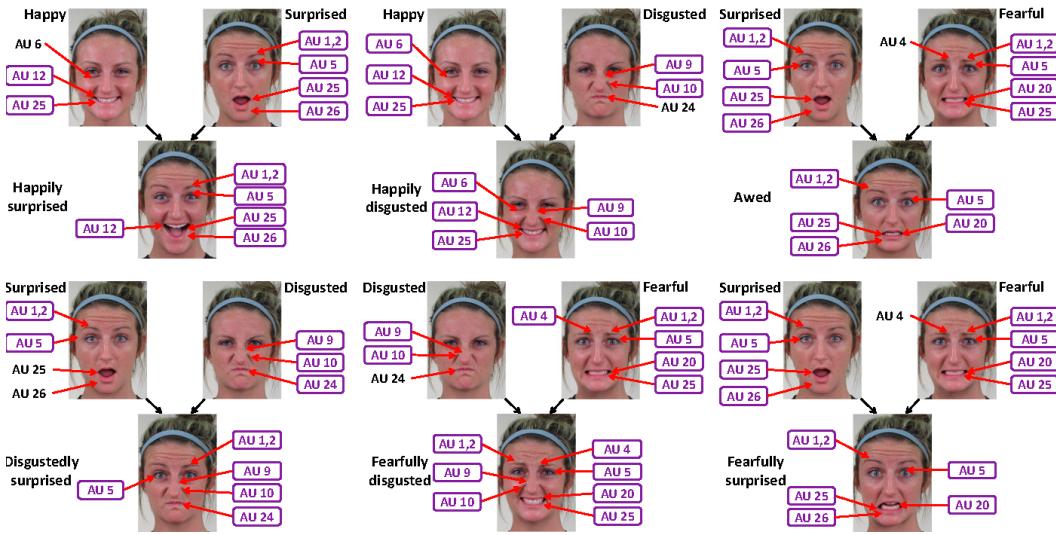


Figure 1.11. AUs of six compound facial expressions of emotion. The AUs of the basic emotions are combined to produce the compound category [29].

Substantial work on AU classification and intensity estimation has been done in [12] by Baltrušaitis et al. while developing the OpenFace [14] system.

Baltrušaitis et al. developed a real-time Facial Action Unit occurrence and intensity detection system based on appearance and geometric features, using Histogram of Oriented Gradients, Shape Parameters and Landmark Locations. They achieved good results by using a median based feature normalization technique so that they could account for person specific neutral expressions.

In [42] Hao et al. explore the dependencies between AUs, and propose a new AU recognition method consisting of a three layer Bayesian network where the top two layers are latent regression Bayesian networks (LRBN).

LRBN are graphical models consisting of a visible layer representing the ground truth for AUs and a latent one. LRBN is able to capture the dependencies between AUs. They tested this system on the CK+ [57], SEMAINE [62] and BP4D [110] database, demonstrating that their approach can accurately capture AU relationships.

In [23] the authors try to model three fundamental aspects of automated AU detection: spatial representation, temporal modeling, and AU correlation. They

proposed an approach using a hybrid network architecture.

Spatial representation is extrapolated by using a Convolutional Neural Network (CNN). For temporal dependencies Long Short-Term Memory Neural Networks (LSTM) are stacked on top of the CNN.

The output of LSTMs and CNNs are then fused together to predict 12 AUs on a frame to frame basis. This network was then tested on the GFT and BP4D [110] dataset, gaining improvements over standard CNN.

De la Torre et al. [21] [22] tackle this problem by personalizing a generic classifier without requiring additional labels for the test subject (unsupervised). They use a transductive learning method, referred as Selective Transfer Machine (STM), that personalizes a generic classifier by attenuating people specific biases. This is done by learning a classifier and re-weighting the training samples most relevant to the subject.

The performance of STM were compared to generic classifiers and cross-domain learning methods and evaluated on CK+ [57], GEMEP-FERA, RUFACS and GFT dataset, and STM is shown to outperform the generic classifiers on all datasets.

Although there has been a lot of progress during the years, there is still a lot of work to be done regarding AUs, since many approaches to automatic facial AU detection fail to account for individual difference in morphology and behavior for the target person. This is a hard problem because it's difficult for classifiers to generalize to very different subjects, and training a person-specific classifier is neither feasible (very low training data) nor there is particular theoretical interest in the research community.

1.4 Contributions and Outline

1.4.1 Contribution

The human face holds many clues to recognize emotions and to discern truth and lies. Human emotions can be recognized mainly by the expressions we make when a situation arises, such as being scared, happy, confused angry and so on. Those expressions are represented by movements of the muscles of the face. Building on this idea, the objective of this thesis is to classify whether a person is lying or telling the truth based on the action units (muscular movements of the face) performed by the subjects and extracted from a video.

It is important to note that all the videos are taken from high stakes situations, specifically court trials, so that all the liars and truth tellers have a big incentive in lying or telling the truth, since the validity of many studies is questionable based on the fact that the emotions analyzed are posed (non spontaneous).

As far as we know this is the first time a study uses action units extracted from a video and processed through an SVM to classify between truth and lies in high stakes situations.

1.4.2 Outline

This thesis is composed by the following chapters:

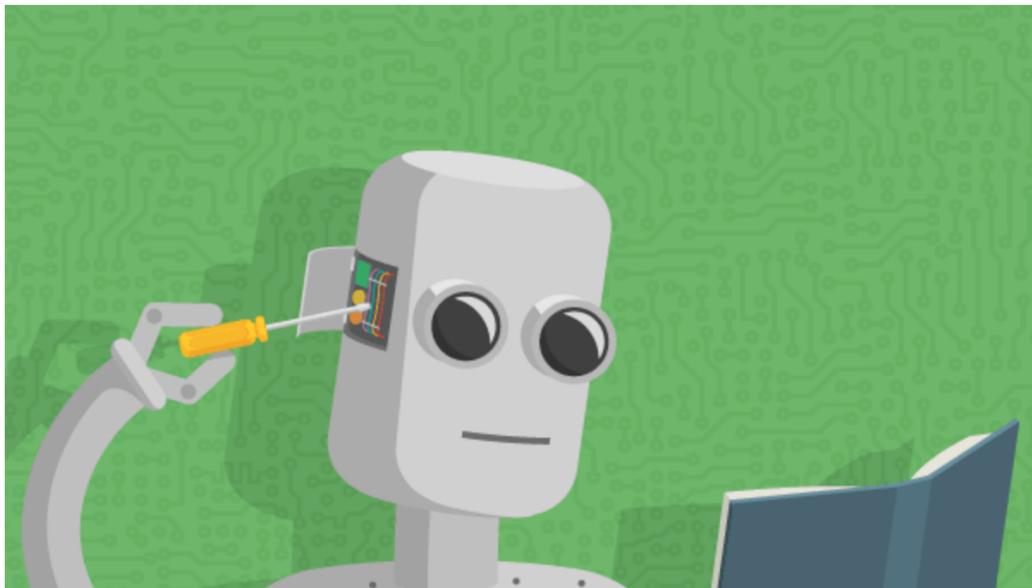
- **Chapter 2** is about Machine Learning. We will talk about classification, regression, supervised learning and some techniques such as Random Forest and SVM.
- **Chapter 3** will describe the architecture of our work, explaining landmark detection, action unit intensity and presence estimation, and the process to discriminate between deception and honesty.
- In **chapter 4** we discuss the dataset and explain the experimental procedure of training and testing a binary classifier for discerning truth and lies. We then proceed discuss the results of the work (par. 4.4)
- In **chapter 5** we draw the conclusions, and future work is presented.

Chapter 2

Machine Learning

In this chapter we will give a brief introduction of Machine Learning (2.1), explaining Classification (2.1.3), Regression (2.1.4) and Supervised learning (2.1.1), proceeding then to explain different techniques such as Random Forests (2.3) and SVM (2.4), which is the most important technique used for this thesis.

2.1 Introduction



Machine Learning is a subfield of Artificial Intelligence that uses statistical techniques to provide computers with the ability to progressively improve performance on different tasks using data, while not being explicitly programmed [24].

The applications for machine learning are huge and diverse, and range from character recognition to email filtering, with lots of applications in computer vision, such as image recognition and classification.

Machine learning also focuses on making prediction on data, by utilizing techniques taken from statistics and mathematical optimization. This has many applications, ranging from health care by predicting risk factors for diseases or gaining insights for prevention, to sports or politics where actual results of games or elections can be predicted.

The field of machine learning is subdivided in two broad categories (Fig. 2.1), supervised learning (par. 2.1.1) and unsupervised learning (par. 2.1.2) based on whether the data is labeled or not.

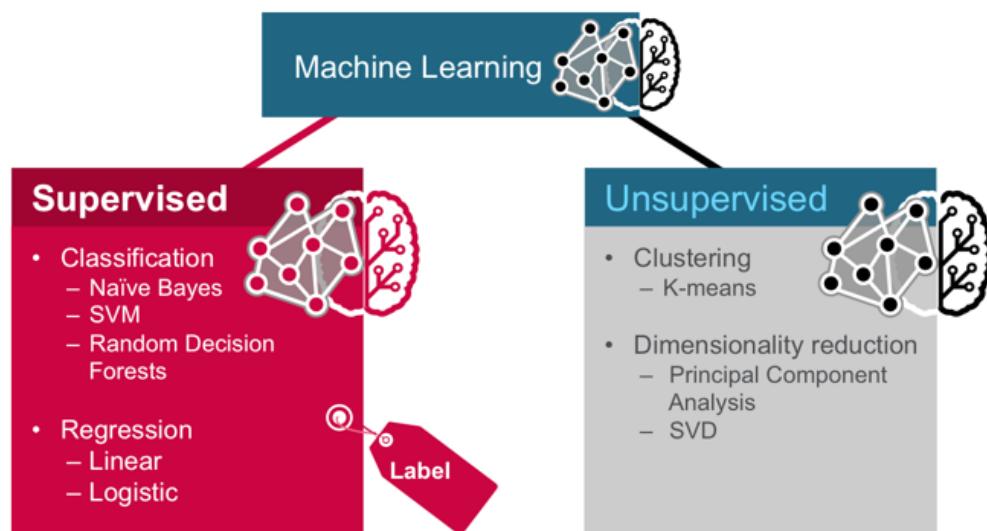


Figure 2.1. Machine Learning Subfields [61]

Another subdivision of the field of machine learning is based on the kind of output that the users wish to obtain, and the categories of said output are:

- **Classification** (par. 2.1.3) for predicting a discrete class label for a new example data.
- **Regression** (par. 2.1.4) to predict a continuous quantity.
- **Clustering** to group similar but unlabeled data.

2.1.1 Supervised Learning

When using supervised learning we want to learn a function that maps an input to an output, based on example input-output pairs [89].

This means that we need labeled training data, consisting of a vector of input objects and an output value. The objective is to correctly classify the new data based on the previously analyzed training pairs.

The general steps to solve a supervised learning problem are (fig. 2.2):

1. Understand what kind of training example to use.
2. Gathering a training set.
3. Model the training set to be fed as input to the algorithm by choosing which features to use and how to represent the data.
4. Choose what kind of algorithm can best train the model.
5. Run the algorithm and evaluate the resulting accuracy on the test set.

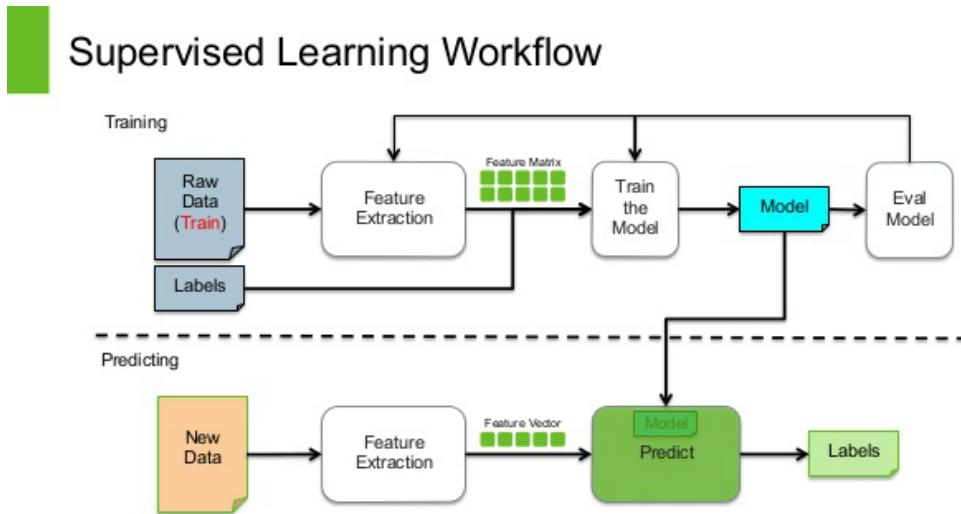


Figure 2.2. Supervised Learning Workflow [59]

There are important considerations to make when using a supervised learning approach:

- **Bias-Variance Tradeoff [81]:**

Suppose we have many different (but equally good) training sets: an algorithm is biased for input x if, when trained on each of these data sets, it is consistently inaccurate at predicting the correct output for x . A high bias indicates that the data points tend to be very close to the mean, and to each other.

A learning algorithm has high variance for a particular input x if it predicts output values that are correct on average but inconsistent when trained on different training sets.

Basically, a high variance indicates that the data points are very spread out from the mean, and from one another (Fig. 2.3).

The prediction error of a classifier is related to the sum of bias and variance, so generally there is a tradeoff between them. Low bias means that it fits the new data well, but if the bias is too low it will fit each training set differently and so result in high variance.

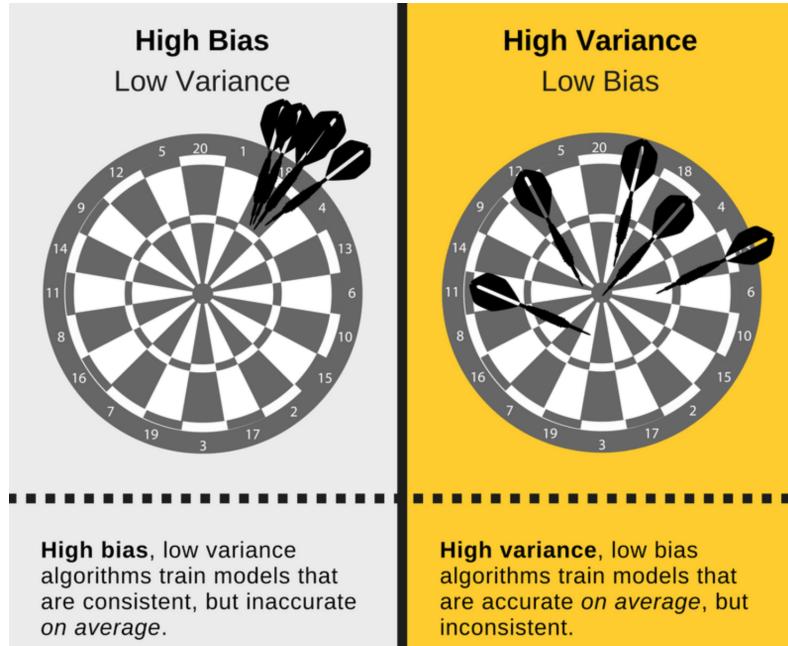


Figure 2.3. Example of Bias-Variance Tradeoff [17]

- **Dimensionality of Input** [99]:

When the input feature vectors are dimensionally very big there could be problems in learning the function, even if not all features contribute significantly to it. This happens because the data depends on too many variables and this could cause high variance.

To avoid this, it is important to reduce the number of features through manual removal or using feature selection algorithms such as Principal Component Analysis (PCA). This usually improves the accuracy of the classifier.

- **Overfitting and Underfitting** [19]:

Overfitting happens when the algorithm adapts too much on the training data and is no longer able to make accurate predictions on the test data (Fig. 2.4). This usually happens when there is an excessive number of parameters than what can be justified by the data [33].

Underfitting is the opposite: a model is not able to correctly capture the structure of the data, for example when fitting non linear data with a linear model.

A common way to avoid overfitting is to resample the data using different techniques, commonly k-fold cross validation or leave-one-out. Other methods include feature removal, early stopping and regularization.

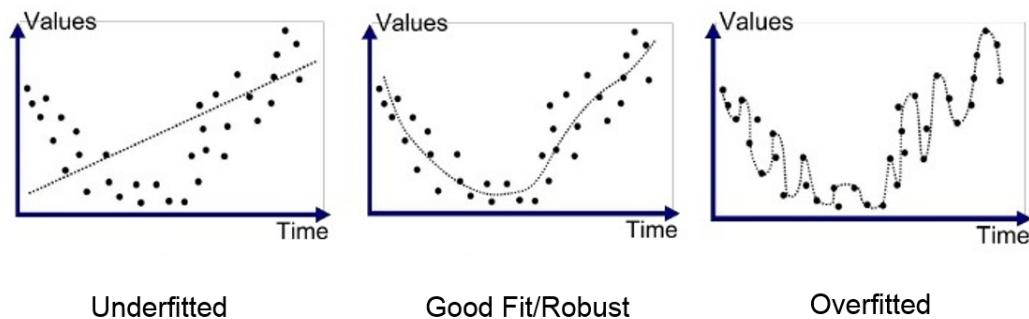


Figure 2.4. Example of overfitting and underfitting [15]

There are many algorithms used to perform supervised learning tasks, the most commonly used are:

- Linear Regression (2.1.4)
- Logistic Regression (2.2)
- Naive Bayes
- Linear Discriminant Analysis
- Decision Trees
- k-Nearest Neighbor
- Neural Networks
- Support Vector Machines (2.4)

2.1.2 Unsupervised Learning

Unsupervised learning is the subfield of Machine Learning tasked with inferring a function from the analysis of unlabeled (not classified) data. Being unclassified there is also a difficulty in evaluating the accuracy of the model.

Usually items are grouped by some measure of similarity, like for example in k-means clustering (Fig. 2.5).

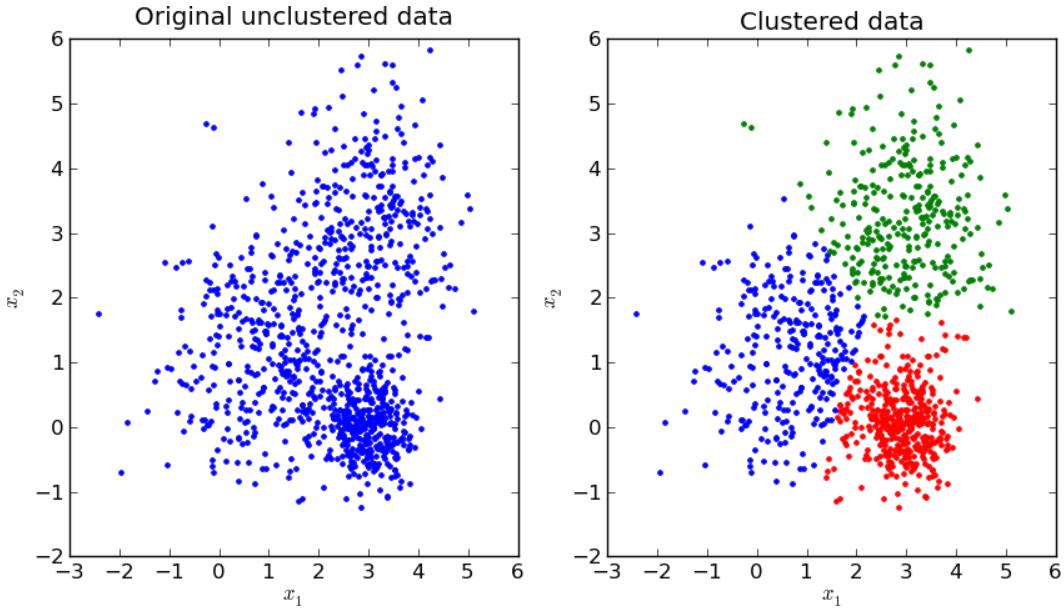


Figure 2.5. Example of Unclustered and Clustered data [48]

These are the most widely used unsupervised learning algorithms:

- Clustering
 - k-means
 - mixture models
 - hierarchical clustering
- Anomaly detection
- Neural Networks
 - Autoencoders
 - Deep Belief Nets
 - Hebbian Learning
 - Generative Adversarial Networks
 - Self-organizing map
- Expectation–maximization algorithm (EM)
- Method of moments
- Blind signal separation techniques
 - Principal component analysis,
 - Independent component analysis,
 - Non-negative matrix factorization
 - Singular value decomposition.

2.1.3 Classification

In machine learning, classification is the problem of identifying in which of a set of categories a new observation belongs, based on a training set of data containing observations whose category is known in advance. A common example could be classifying an email as spam or not [107].

Classification is considered an instance of supervised learning, based on instances where a training set is available. As for unsupervised learning, classification would be clustering since it groups data into categories based on similarity, but without knowing the label of the data.

The observations, called features or explanatory variables, take different types based on the value. They can be categorical, numerical or ordinal, or be compared by similarity between previous observations using some kind of distance function.

The observations representing the categories to be predicted are called explanatory variables (or regressors, or independent variables).

The classifier is the algorithm that implements the classification, or a function that maps input data to a category.

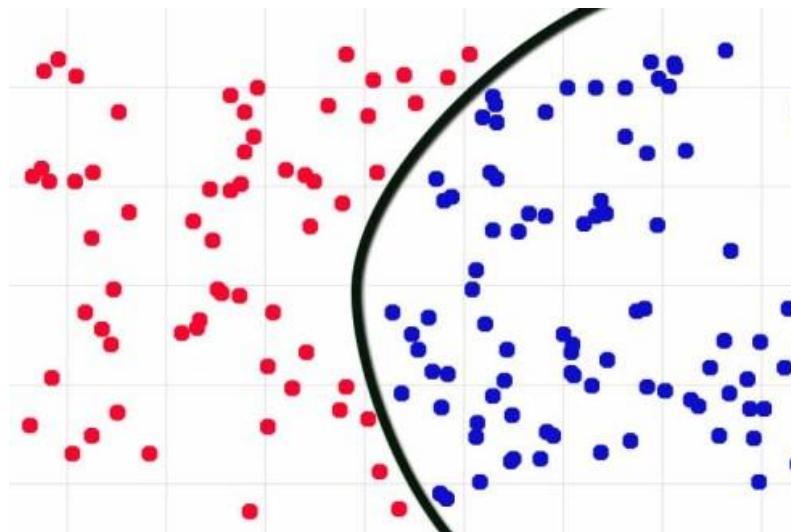


Figure 2.6. Example of data division by a classification algorithm

Many classification algorithms, such as SVM (2.4), logistic regression 2.2, LDA (Linear Discriminant Analysis) or perceptron, can be described using a linear function assigning a score to each category c , by doing the dot product of the feature vector of an instance with a vector of weights, thus combining them. The predicted category will be the one with the highest score. This function is called linear predictor function and has this general formula:

$$\text{score}(X_i, c) = \beta_c \cdot X_i \quad (2.1)$$

where X_i is the feature vector of instance i and β_c is the vector of weights of category c . This kind of algorithms are known as linear classifiers.

2.1.4 Regression Analysis

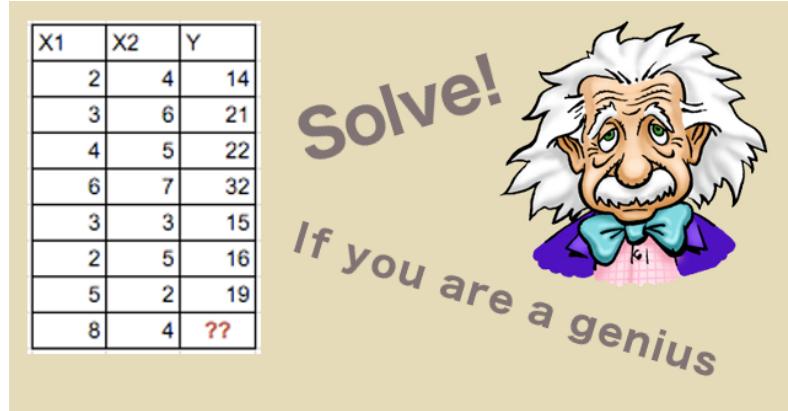


Figure 2.7. Regression Analysis can solve this! [10]

Regression analysis is used in statistics and machine learning to estimate the relationships among variables.

More specifically, it focuses on the relationship between one dependent variable and more independent variables (also called predictors), and generally is used to estimate the average value of the dependent variable when the independent variables are fixed, by calculating a regression function.

By understanding the relationship between dependent and independent variables it's also possible to perform predictions on future data based on new inputs, or on changes to the old ones [106].

In regression analysis, it is also of interest to characterize the variation of the dependent variable around the prediction of the regression function using a probability distribution.

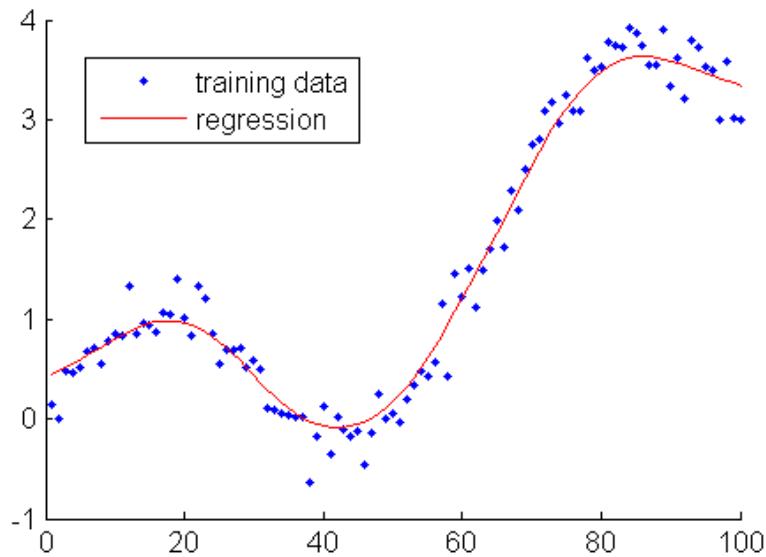


Figure 2.8. Example of Regression Analysis

Regression Model

A general regression model uses the following variables and parameters:

- The unknown parameters, β , that represents either a scalar or a vector.
- The independent variables, X .
- The dependent variable, Y .

A regression model relates Y to a function of X and β [106].

$$Y \approx f(X, \beta) \quad (2.2)$$

This is usually formalized as:

$$E(Y|X) = f(X, \beta) \quad (2.3)$$

If β is of length k and the number of observed data points is enough ($N > k$), then it's possible to estimate a unique value for β that best fits the data.

If this is the case, regression analysis provides the means to find a solution for unknown parameters of β to, for example, apply the method of least squares.

To apply regression the data must abide by some assumptions, generally [101]:

- The sample is representative of the population for the inference prediction.
- The error is a random variable with a mean of zero.
- The independent variables are measured with no error.
- The independent variables (predictors) are linearly independent.
- The errors are uncorrelated.
- The variance of the error is constant across observations.

2.1.5 Linear Regression

Linear regression is the most basic case of Regression Analysis, and it is a useful tool for predicting a quantitative response. Also, most of the newer approaches to regression are often a generalization or extension of linear regression.

Linear regression is a linear approach to modeling the relationship between a dependent variable and one or more independent variables. The simplest case of linear regression (LR) is when there is only one independent variable, and is called simple linear regression [101] and has this form:

$$y_i \approx \beta_0 + \beta_1 x_i + \varepsilon_i \quad (2.4)$$

where ε_i is the error for the i th observation. The error is a catch-all for what we miss with this simple model, since it's very probable that the true relationship is not linear, and that there may be other variables that influence y [37].

The aim of LR is to estimate the β coefficients to make predictions. To do so, we utilize of the training dataset to produce estimates $\widehat{\beta}_0$ and $\widehat{\beta}_1$ for the model coefficients. We can then make predictions by computing:

$$\widehat{y}_i \approx \widehat{\beta}_0 + \widehat{\beta}_1 x_i \quad (2.5)$$

$e_i = y_i - \widehat{y}_i$ is the difference between the true value and the prediction for an observation, and it is called residual (Fig. 2.9).

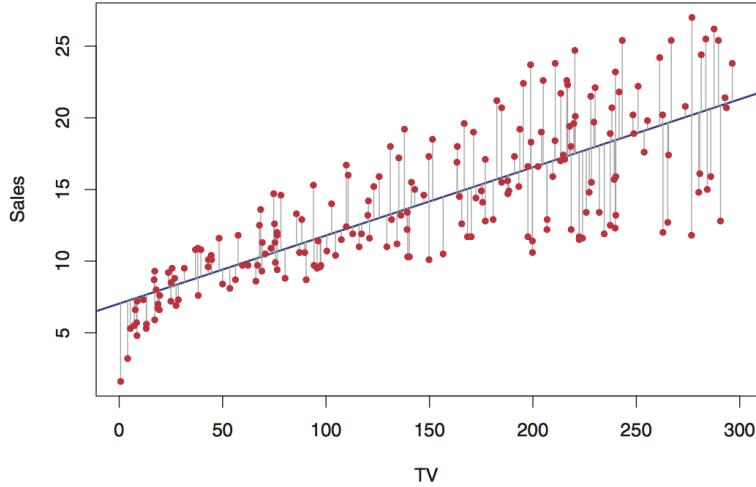


Figure 2.9. The fit is found by minimizing the sum of squared errors. Each gray line segment represents an error, and the prediction makes a compromise by averaging their square [37]

The most common method for estimation is called least squares. This method obtains parameter estimates that minimize the sum of squared residuals (SSR):

$$SSR = \sum_{i=1}^n e_i^2 \quad (2.6)$$

The minimizers are

$$\widehat{\beta}_1 = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2} \quad (2.7)$$

$$\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \bar{x} \quad (2.8)$$

where \bar{x} and \bar{y} are the mean of x and y .

If we assume that the error term has constant variance, the estimate of the variance of the error is given by:

$$\widehat{\sigma}_\varepsilon^2 = \frac{SSR}{n - 2} \quad (2.9)$$

And is called mean square error (MSE) of the regression. The denominator is the sample size reduced by the number of model parameters estimated from the same data, $(n - p)$ for p regressors or $(n - p - 1)$ if an intercept is used [88]. In the case of simple linear regression $p = 1$, so the denominator is $n - 2$.

We can then estimate the standard errors for the parameters, that tell us the average amount that the estimate differs from the actual value.

$$\widehat{\sigma}_{\beta_1} = \widehat{\sigma}_\varepsilon \sqrt{\frac{1}{\sum(x_i - \bar{x})^2}} \quad (2.10)$$

$$\widehat{\sigma}_{\beta_0} = \widehat{\sigma}_\varepsilon \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2}} \quad (2.11)$$

Standard errors can be used to compute confidence intervals. A 95% confidence interval is defined as a range of values such that with 95% probability, the range will contain the true unknown value of the parameter.

The general multiple regression model, where there are p independent variables, follows this formula:

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} + \varepsilon_i \quad (2.12)$$

where x_{ij} is the i-th observation on the j-th independent variable.

2.2 Logistic Regression

In regression analysis, logistic regression is a method for estimating the parameters of a logistic model. A logistic model is one where the log-odds of the probability of an event is a linear combination of independent or predictor variables.

The two possible dependent variable's values are often labeled as "0" and "1", "true" or "false", "pass" or "fail" etc, and they represent binary outcomes (Fig. 2.10).

It's possible to generalize the binary logistic regression model to more than two levels (outcomes) of the dependent variable [102].

The binary logistic model works by estimating the probability of a binary response based on one or more predictor variables.

By using logistic regression is possible to say how much that the presence of a risk factor increases the odds of a given outcome.

The model itself simply calculates the probability of output in terms of input, and it is not specifically a classifier, but it can be used to classify data by choosing a threshold value, for example if the probability is $\geq 50\%$ the class is "1" otherwise is "0".

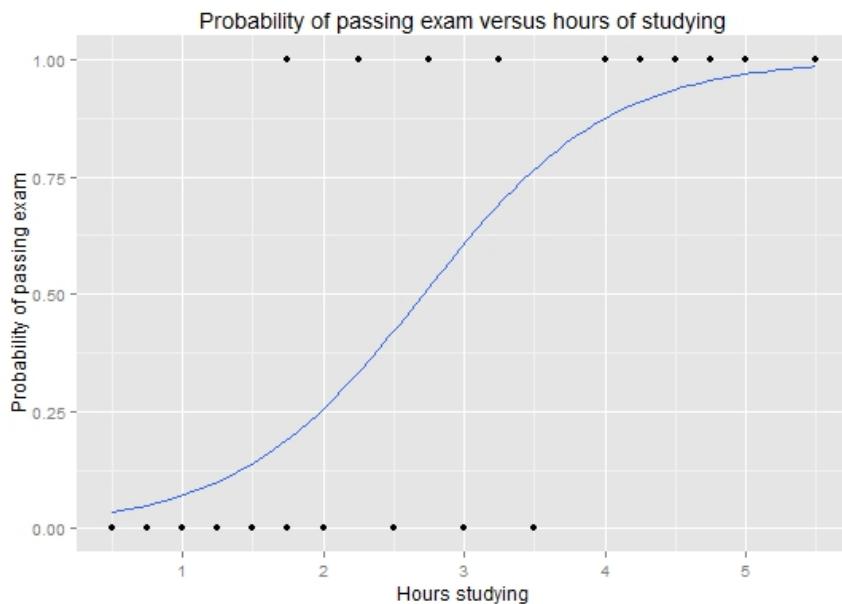


Figure 2.10. Graph of a logistic regression curve showing probability of passing an exam versus hours studying [102]

Logistic Function

The logistic function to estimate the probabilities is a sigmoid function (exemplified in fig. 2.11).

A sigmoid function takes an input from the real numbers and outputs a value between 0 and 1.

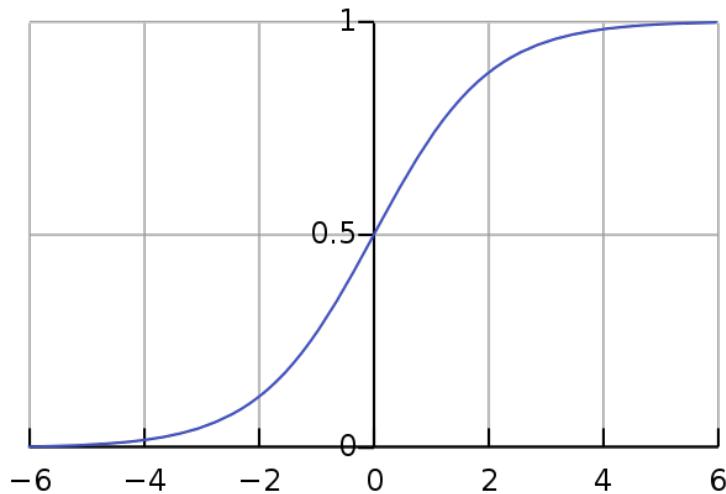


Figure 2.11. Example of a Sigmoid shaped function

We want to model the relationship between $p(X) = \Pr(Y = 1|X)$ and X . To model this relationship we use the a logistic function:

$$p(X) = \frac{e^{\beta_0} + e^{\beta_1 X}}{1 + e^{\beta_0} + e^{\beta_1 X}} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}} \quad (2.13)$$

Equation 2.13 can be manipulated into:

$$\frac{p(X)}{1 - p(X)} = e^{\beta_0} + e^{\beta_1 X} \quad (2.14)$$

$\frac{p(X)}{1 - p(X)}$ is called odds, and can take any value between 0 and ∞ .

By doing the logarithm of 2.14 we get:

$$\log \frac{p(X)}{1 - p(X)} = \beta_0 + \beta_1 X \quad (2.15)$$

The coefficients β_0 and β_1 in equation 2.13 are unknown, and must be estimated based on the available training data.

To fit the model of equation 2.15, the maximum likelihood method is commonly used.

Intuitively we want to estimates β_0 and β_1 such that the predicted probability $\hat{p}(x_i)$ for a specific case corresponds as closely as possible to the observed class for that case.

This can be achieved by using the likelihood function:

$$\ell(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i':y_{i'}=0} (1 - p(x_{i'})) \quad (2.16)$$

The estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ are chosen to maximize this likelihood function.

Maximum likelihood is a very general approach that is used to fit many of the non-linear models. For example in linear regression (par. 2.1.5), the least squares approach is a special case of maximum likelihood.

Once the coefficient have been estimated is pretty simple to compute the probability:

$$\hat{p}(X) = \frac{e^{\hat{\beta}_0} + e^{\hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0} + e^{\hat{\beta}_1 X}} \quad (2.17)$$

2.3 Random Forest

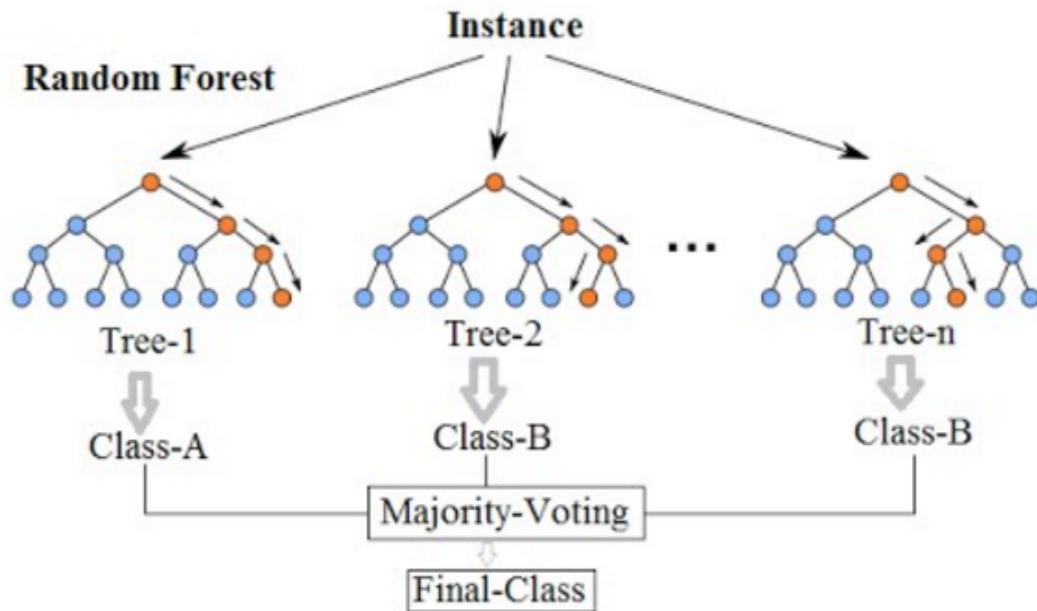


Figure 2.12. Example of Random Forest Classification [51]

Random forests are an ensemble learning method used mainly for classification and regression. They work by building many decision trees while in the training phase, and then outputting the class that is the mode of the classes for classification, or the mean prediction of the individual trees for regression [105][46].

Random Forest (RF) is preferred to decision trees because decision trees suffer from high variance. By using tree bagging, which is a commonly used technique to reduce variance, we can correct decision trees' habit of overfitting to their training set [44].

Decision Trees

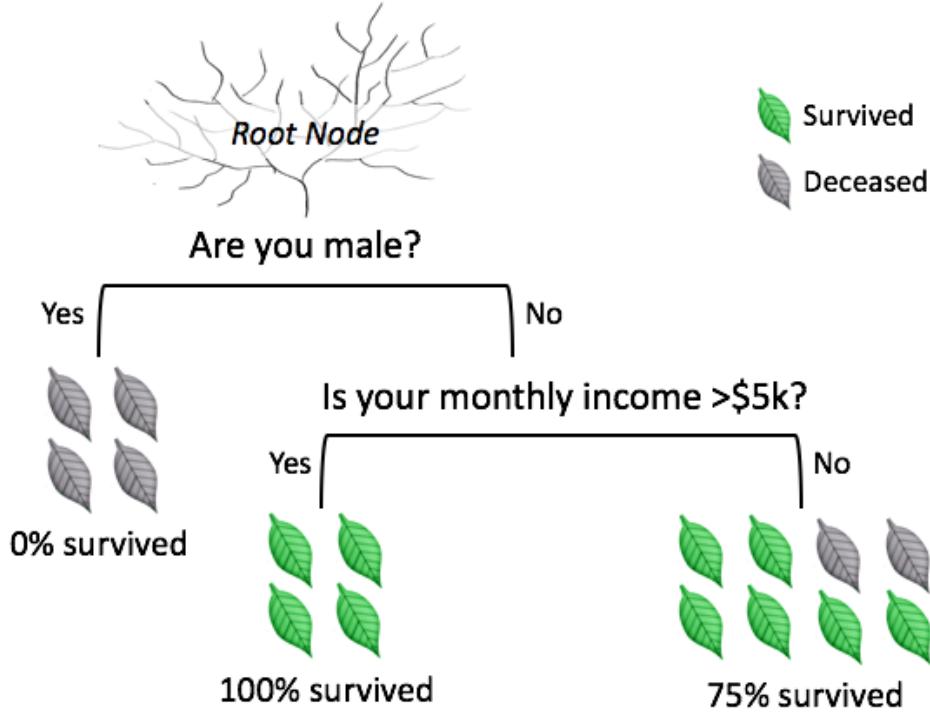


Figure 2.13. Example of a Decision Tree [67]

Decision trees are a widespread method for many different machine learning tasks and can be used for either classification or regression.

The aim of decision trees is to build a model that predicts the value of a target variable based on different input variables.

Each node of the tree corresponds to one of the input variables. Each leaf represents a value of the target variable, given the values of the input variables represented by the path from the root to the leaf (fig 2.13) [98].

The data usually has the following structure:

$$(x, Y) = (x_1, x_2, \dots, x_n, Y) \quad (2.18)$$

Where Y is the dependent variable that we are trying to classify and x is the feature vector.

A decision tree is able to represent a classification problem: The arcs of a tree's node labeled with input features are labeled with the possible values for that feature, and they can lead to two options:

- A subordinate decision node, with different input features.
- A leaf labeled with the resulting class (or a probability).

A tree can be learned by dividing the training set into subsets based on the values of the attributes. This is done recursively until a node has all the same value of the target variable, or when splitting is no longer aiding the prediction's result.

Decision trees, however, are often inaccurate. Specifically, by growing trees very deep, they tend to learn highly irregular patterns by overfitting to the training set (they normally have low bias but very high variance) [44].

Random forests are a way of averaging multiple deep decision trees, trained on different parts of the same training set, with the goal of reducing the variance. This comes at the expense of a small increase in the bias and some loss of interpretability, but generally greatly boosts the performance in the final model [44].

Tree Bagging

Given a set of n independent observations X_1, \dots, X_n , each with variance σ^2 , the variance of the mean \bar{X} of the observations is given by σ^2/n .

This means that averaging a set of observations reduces the variance. So to reduce the variance we can take many training sets from the population, build a prediction model using each training set, and average the resulting predictions.

Basically we could calculate $\hat{f}^1(x), \hat{f}^2(x), \dots, \hat{f}^B(x)$ using B separate training sets, and then average them in order to obtain a single model with lower variance [37].

This solution is not practical since we normally lack more than one training set. Instead we utilize *bagging*: we can take repeated samples from the training set in order to generate B different bootstrapped training sets. We then train our method on the b th bootstrapped training set, and then average the resulting predictions.

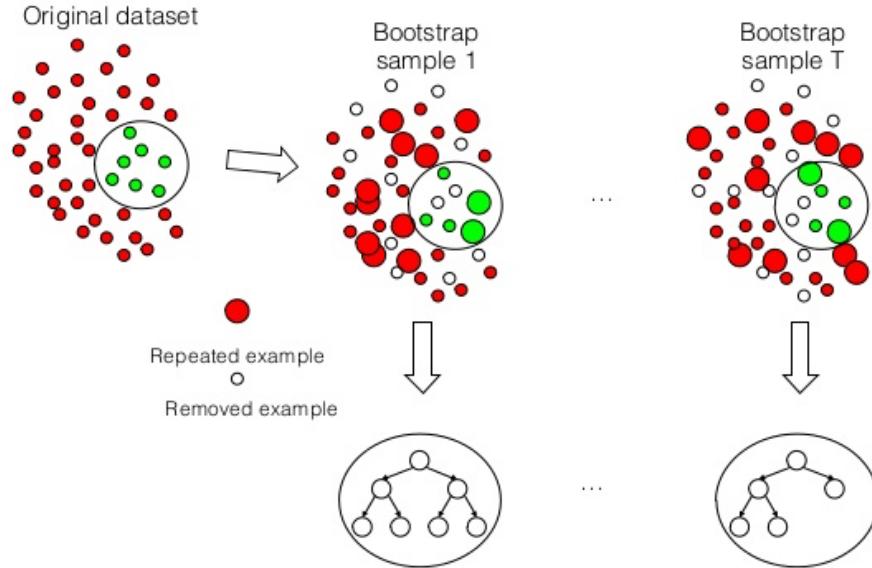


Figure 2.14. Bagging Example [1]

To apply bagging to regression trees specifically, we simply construct B regression trees using B bootstrapped training sets, and average the resulting predictions. These trees are grown deep, and are not pruned. Hence each individual tree has high variance, but low bias. Averaging these B trees reduces the variance [37].

If we are using Random Forests for classification (like in this thesis), the simplest approach is the following:

Given a test observation, we record the class predicted by each of the B trees, and take a majority vote: the final prediction is the most frequent class among the B predictions.

More formally, given a training set $X = x_1, \dots, x_n$ with responses $Y = y_1, \dots, y_n$, bagging B times selects a random sample with replacement of the training set and fits trees to these samples [105]:

For $b = 1, \dots, B$:

1. Sample, with replacement, n training examples from X, Y ; call these X_b, Y_b .
2. Train a classification or regression tree f_b on X_b, Y_b .

After training, predictions for unseen samples x' can be made by averaging the predictions from all the individual regression trees on x' :

$$\hat{f} = \frac{1}{B} \sum_{b=1}^B f_b(x') \quad (2.19)$$

or by using the majority vote in case of classification.

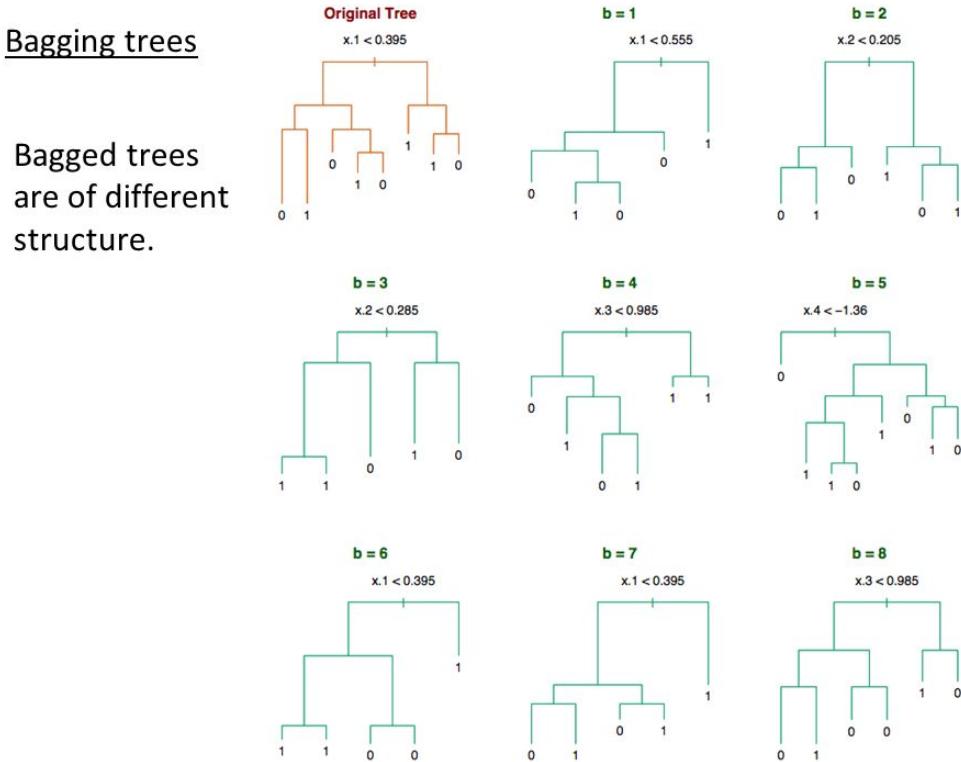


Figure 2.15. Bagging leads to different decision trees [2]

Additionally, an estimate of the uncertainty of the prediction can be made as the standard deviation of the predictions from all the individual regression trees on x' :

$$\sigma = \sqrt{\frac{\sum_{b=1}^B (f_b(x') - \hat{f})^2}{B - 1}} \quad (2.20)$$

The number of samples/trees, B , is a free parameter. Typically, a few hundred to several thousand trees are used, depending on the size and nature of the training set. An optimal number of trees B can be found using cross-validation, or by observing the out-of-bag error (2.3): the mean prediction error on each training sample x_i , using only the trees that did not have x_i in their bootstrap sample.

The training and test error tend to level off after some number of trees have been fit [37].

From Bagging to Random Forests

Random Forests provide an improvement over bagged decision trees by decorrelating the trees. As in bagging, the decision trees are built on a bootstrapped training set. The difference is that every time a split is considered, a random sample of m predictors is chosen as split candidates from the total p predictors. The split can only use one of those m predictors. \sqrt{m} predictors are taken each split, typically with $m \equiv \sqrt{p}$ [37].

Practically, at each split in the tree, the algorithm can only consider a subset of the predictors (Fig. 2.16). This works because if there is a very strong predictor, that predictor would be used in the majority of the trees, and then all the bagged trees would look alike and would be highly correlated, thus leading to a small reduction in variance.

By forcing each split to consider only m predictors, on average $(p - m)/p$ of splits will not take into account a strong predictor. This, in turn, will decorrelate the trees.

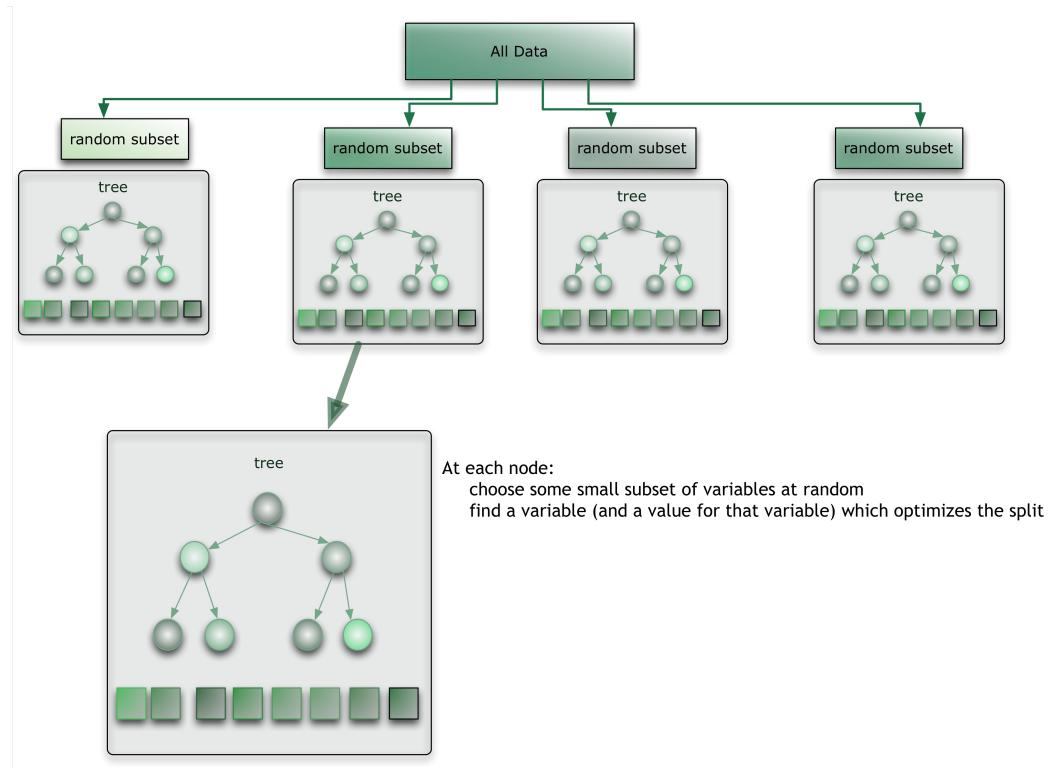


Figure 2.16. Random Forest overview

Out of bag Error

There is a very direct way to estimate the test error of a bagged model: on average, we use $2/3$ of the observations to fit a bagged tree. The remaining $1/3$ that are not utilized are called the out-of-bag (OOB) observations.

We can predict the response for the i th observation using each of the trees in which that observation was OOB. This will generate approximately $B/3$ predictions for the i th observation. We can then average these predictions (or use majority voting for classification).

We can repeat the process for all the n observations and get an average classification error.

The resulting OOB error is a valid estimate of the test error for the bagged model, since the response for each observation is predicted using only the trees that were not fit using that observation [37].

Variable Importance

Bagging has the advantage of giving increased accuracy over using a single tree, at the cost of interpretability of the data. Using bagging we can no longer have easily interpreted diagrams.

With Random Forest we can instead rank the importance of each predictor: given a dataset $D_n = \{X_i, Y_i\}_{i=1}^n$ we fit a RF to the data. While fitting the data, the out-of-bag error is averaged over the forest.

To get the importance of the j th feature, we permute it's values over the training data and recompute the OOB error.

The importance is given by doing the average of the difference between the OOB error before and after this permutation. Then the score is normalized by the standard deviation of the differences.

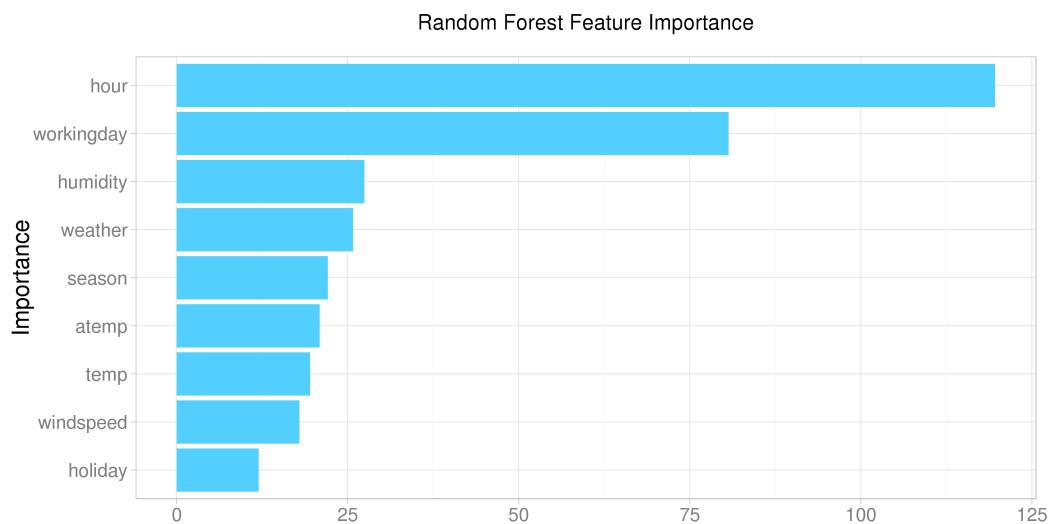


Figure 2.17. Example of variable importance found with Random Forest [40].

2.4 SVM

Support Vector Machines (SVM) are a supervised machine learning algorithm used for both classification and regression.

The main idea is to find the optimal hyperplane for linearly separable data, and then extend this idea to data that are not linearly separable by mapping this data in a new space using a kernel function.

2.4.1 Hyperplane

What is a hyperplane? In p -dimension, a hyperplane is a flat affine subspace of $p - 1$ dimension. For example, in two dimensions it's a line, in three dimensions it's a plane (Fig. 2.18), and so on.

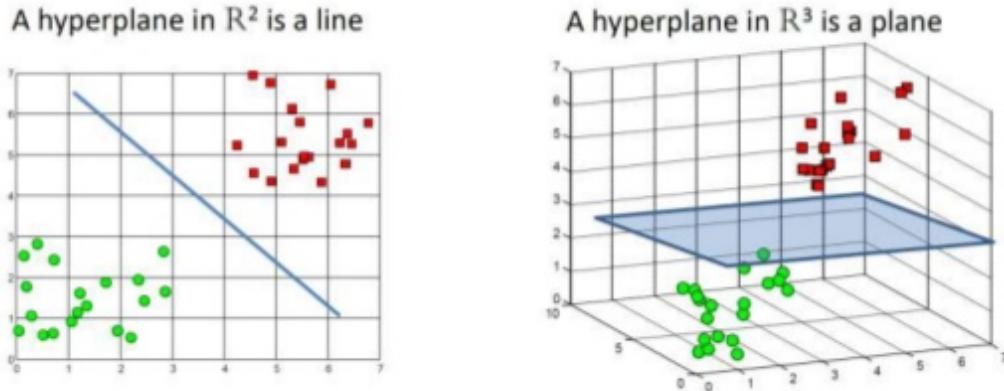


Figure 2.18. Two and three dimensional space hyperplanes [84]

The definition of an hyperplane for two dimensions is:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0 \quad (2.21)$$

Where any $X = (X_1, X_2)^T$ for which equation 2.21 holds is a point on the hyperplane. And it can be extended to an hyperplane of p -dimensions:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p = 0 \quad (2.22)$$

If X does not satisfy equation 2.22, but

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p > 0 \mid < 0 \quad (2.23)$$

we know that X lies to one side or the other of the hyperplane. So a hyperplane divides a space in two halves [37].

It's possible to build a hyperplane that can separate training data based on their class labels, for example by assigning $\{-1, 1\}$ if the result lies on one side or the other of the hyperplane.

2.4.2 Maximal Margin Classifier

If the data can be separated by a hyperplane, it means that there can be an infinite number of separating hyperplanes.

To build a classifier it's necessary to chose one between those infinite hyperplanes, such that the classification results are as accurate as possible.

The distance between the hyperplane and the nearest data point from either side is known as the margin.

Generally the choice is the *maximal margin hyperplane*, meaning the hyperplane that *maximizes* the margins (has the farthest minimum distance) for the data we are classifying.

In fact the distance of a data point from the hyperplane is a measure of our confidence that the observation was classified correctly [37].

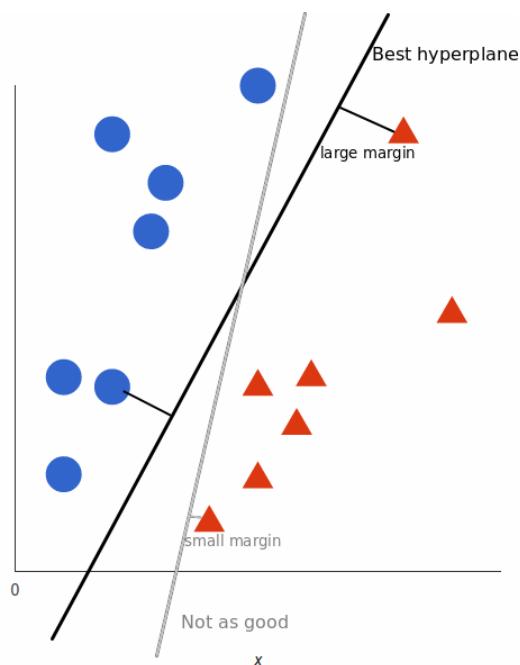


Figure 2.19. Different hyperplanes have different margins [87]

The idea is that a classifier with a large margin on the training set will have a large margin on the test set, so it will classify it correctly [37].

2.4.3 Support Vectors

Support vectors are the data points that lie closest to the hyperplane (Fig. 2.20), they are also the most difficult data points to classify and have direct influence on the location of the hyperplane. In fact, moving the support vectors would change the hyperplane's location.

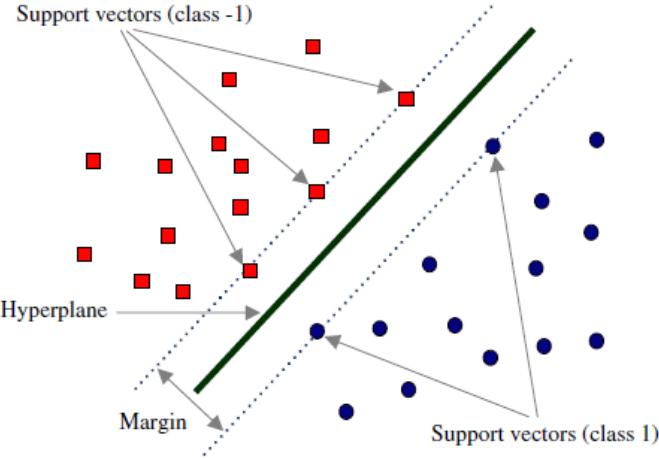


Figure 2.20. Example of Support Vectors

So the objective is to choose a hyperplane with the largest possible margin between it and the support vectors, since the larger is the margin, the lower the generalization error of the classifier [37].

2.4.4 Finding the Maximal Margin Classifier

Finding the maximal margin hyperplane based on a set of n training observations $x_1, \dots, x_n \in \mathbb{R}^p$ and with class labels $y_1, \dots, y_n \in \{1, -1\}$, translates to an optimization problem:

Maximize the margin M :

$$\beta_0, \beta_1, \beta_2, \dots, \beta_p, M \quad (2.24)$$

subject to

$$\sum_{j=1}^p \beta_j^2 = 1 \quad (2.25)$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M \quad \forall i = 1, \dots, n \quad (2.26)$$

Equations 2.25 and 2.26 ensure that each observation is on the correct side of the hyperplane, and at least at distance M from the hyperplane.

So M represents the margin of our hyperplane, and the optimization problem chooses $\beta_0, \beta_1, \beta_2, \beta_3, \dots, \beta_p$ to maximize M and find the maximal margin hyperplane [37].

2.4.5 Support Vector Classifier

Unfortunately this hyperplane does not necessarily exists (fig. 2.21), but we can extend the concept of the Maximal Margin Classifier to find a hyperplane that almost separates the classes by using a *soft* margin. This is actually what an SVM does.

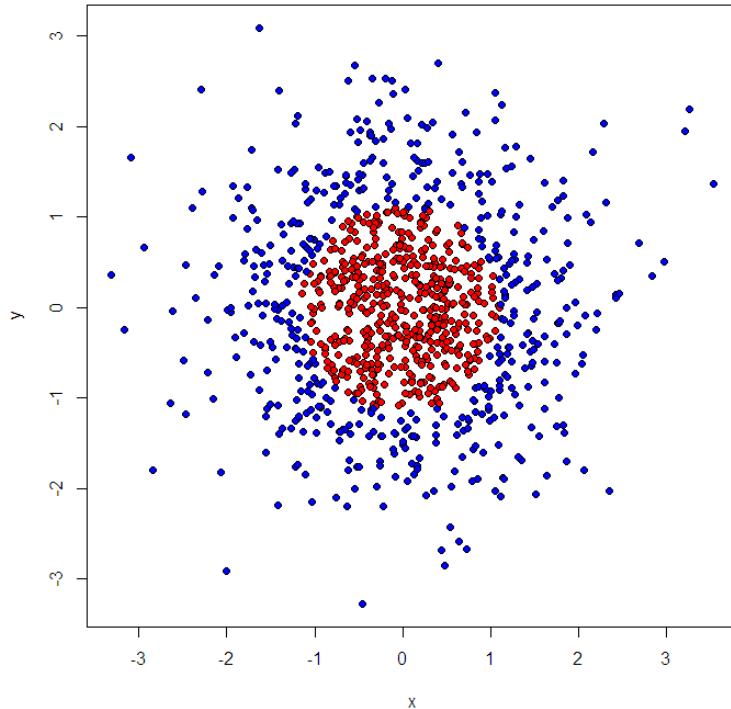


Figure 2.21. Example of Linearly non separable data [78]

Even if a separating hyperplane did exists, we might not want to use it: such a classifier would necessarily classify all of the training data perfectly, and that would lead to sensitivity to changes in data points (fig. 2.22). Also, being so sensitive could mean that it's overfitting the data.

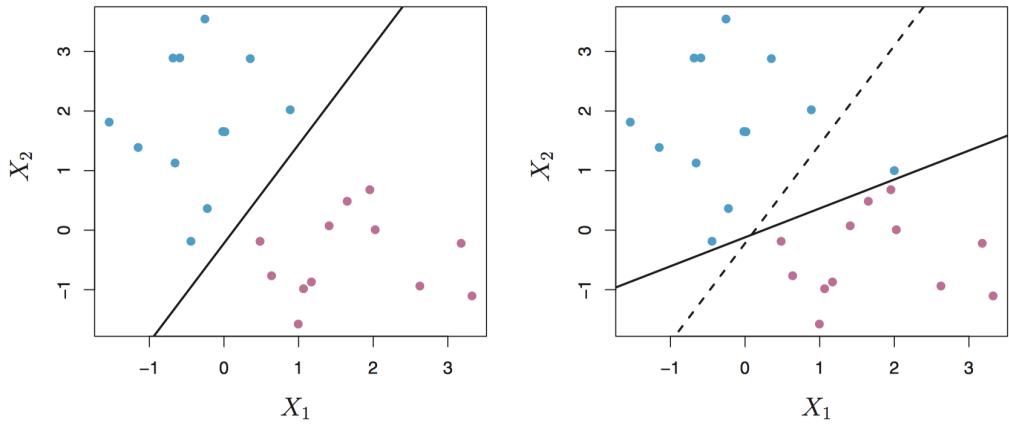


Figure 2.22. Adding a single data point (from left to right) can change the maximal margin hyperplane very significantly [37]

Instead we consider a classifier that does not separate the data perfectly, but is less prone to change for individual data points and can better classify *most* of the observations.

This classifier is called *soft margin classifier* or *support vector classifier* (Fig. 2.23). In a support vector classifier (SVC) a small subset of the observations are on the wrong side of the margin, or even on the wrong side of the separating hyperplane.

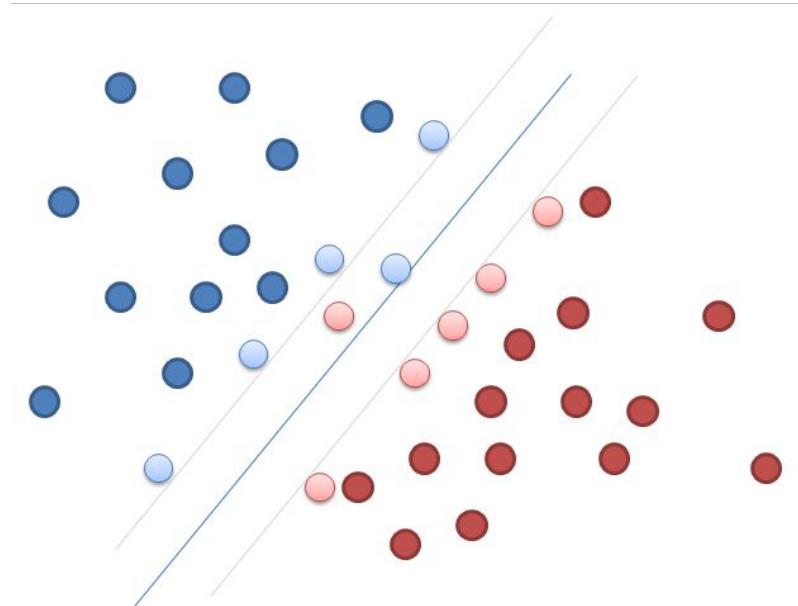


Figure 2.23. Soft margin classifier. Some data points are allowed to be misclassified to find a larger margin [80]

The optimization problem we are solving now is the following:

Maximize the margin M :

$$\beta_0, \beta_1, \beta_2, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n, M \quad (2.27)$$

subject to

$$\sum_{j=1}^p \beta_j^2 = 1 \quad (2.28)$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i), \quad \epsilon_i \geq 0, \quad \sum_{i=1}^n \epsilon_i \leq C \quad (2.29)$$

With $C \geq 0$, and $\epsilon_1, \dots, \epsilon_n$ being slack variables to allow data points to be on the wrong side of the margin.

When we solve equations 2.27-2.29 we can classify x^* by determining where it lies on the hyperplane:

$$f(x^*) = \beta_0 + \beta_1 x_1^* + \beta_2 x_2^* + \dots + \beta_p x_p^* \quad (2.30)$$

Based on the value of ϵ_i we can make some assumptions:

- If $\epsilon_i = 0$ the i th observation is on the correct side of the margin.
- If $\epsilon_i > 0$ the i th observation is on the wrong side of the margin.
- If $\epsilon_i > 1$ the i th observation is on the wrong side of the hyperplane.

C determines the quantity and quality of the violations of the margin and hyperplane that is tolerable by our classifier, and is called *cost*.

- If $C = 0$ then we cannot violate the margin, and $\epsilon_1 = \dots = \epsilon_n = 0$.
- If $C > 0$ then the observations $\leq C$ can be on the wrong side of the hyperplane.

When C is small the margin is small and we have a classifier that is very fit to the data, with low bias and high variance.

When C is big the margin is larger and the classifier is less fit to the data and has high bias and low variance.

The interesting property of this optimization problem is that only the support vectors (observations that are on the margin or that violate the margin) affect the hyperplane, and thus the classifier obtained.

All the observations that lie on the correct side of the margin do not affect the support vector classifier.

2.4.6 Kernels

Sometimes data is not linearly separable (Fig. 2.21). When that happens, Support Vector Machines work by enlarging the feature space using specific function called kernels, utilized to accommodate a non-linear boundary between the classes.

The solution to equations 2.27 - 2.29 involves only the inner products of the observations, and not the observations themselves [37].

We define the inner product between two r -vectors a and b as $\langle a, b \rangle = \sum_{i=1}^r a_i b_i$. So the inner product of two observations x_i, x'_i is given by:

$$\langle x_i, x'_i \rangle = \sum_{j=1}^p x_{ij} x'_{i'j} \quad (2.31)$$

the linear SVC can then be described as:

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i \langle x, x_i \rangle \quad (2.32)$$

To estimate β_0 and $\alpha_1, \dots, \alpha_n$ all that is required are the $\binom{n}{2}$ inner products $\langle x_i, x_{i'} \rangle$ between all pairs of training observations [37].

Every time we compute the inner product (2.31) in equation 2.32 we could replace it with a generalization of this form:

$$K(x_i, x'_i) \quad (2.33)$$

Where K is a kernel function. A kernel function is used to quantify the similarity between two observations.

For example a linear kernel of p features is:

$$K(x_i, x'_i) = \sum_{j=1}^p x_{ij} x'_{i'j} \quad (2.34)$$

or for a polynomial kernel of degree d :

$$K(x_i, x'_i) = (1 + \sum_{j=1}^p x_{ij} x'_{i'j})^d \quad (2.35)$$

or a radial kernel:

$$K(x_i, x'_i) = \exp(-\gamma \sum_{j=1}^p (x_{ij}, x'_{i'j})^2), \quad \gamma \geq 0 \quad (2.36)$$

another example is the sigmoid kernel:

$$K(x_i, x'_i) = \tanh\left(-k \sum_{j=1}^p (x_{ij}, x'_{i'j})^2 + \delta\right) \quad (2.37)$$

We can visualize some of these kernels in fig. 2.24

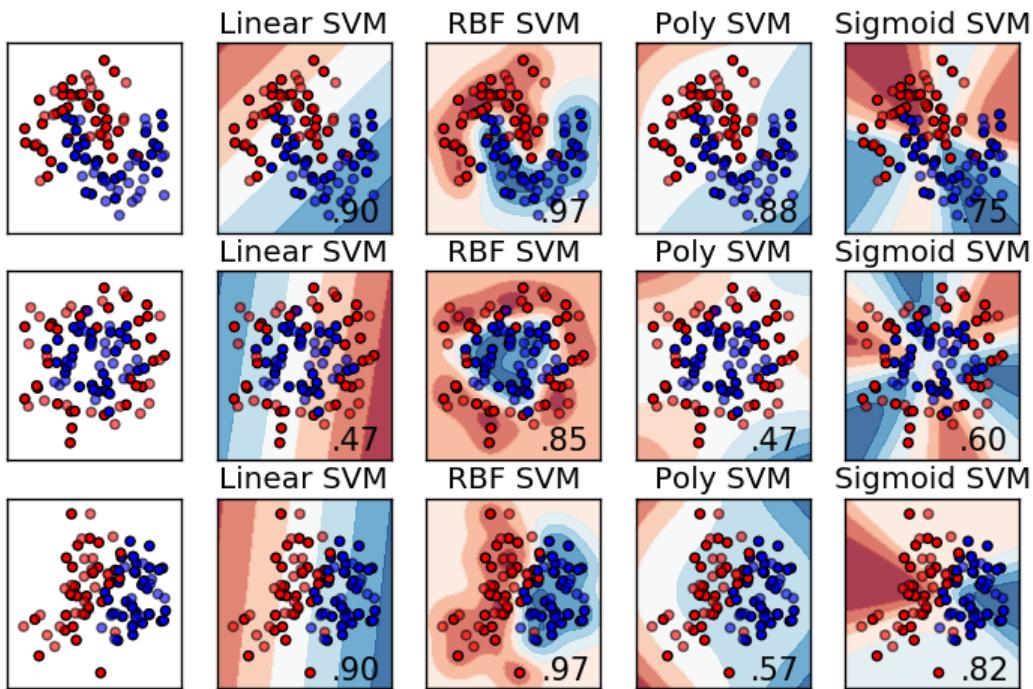


Figure 2.24. Visualization of different kinds of kernels [34].

There are many more kernel functions that can be used based on the shape of the data.

The advantage of using a kernel is first of all computational: by using a kernel we just need to compute $K(x_i, x'_i)$ for all $\binom{n}{2}$ distinct pairs (i, i') .

Graphically, starting with a two dimensional non linearly separable data, this is what happens:

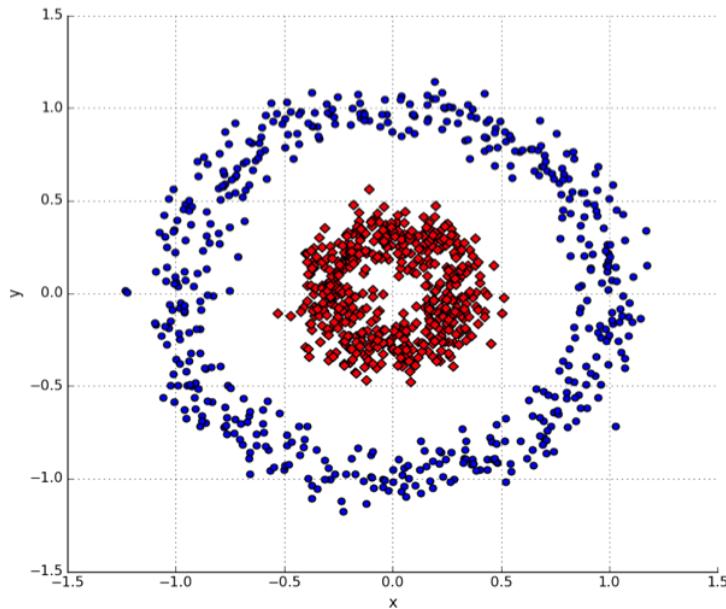


Figure 2.25. Linearly non separable data [92].

We map the data from input space X into a transformed feature space H , by using a non-linear function $\phi : X \rightarrow H$

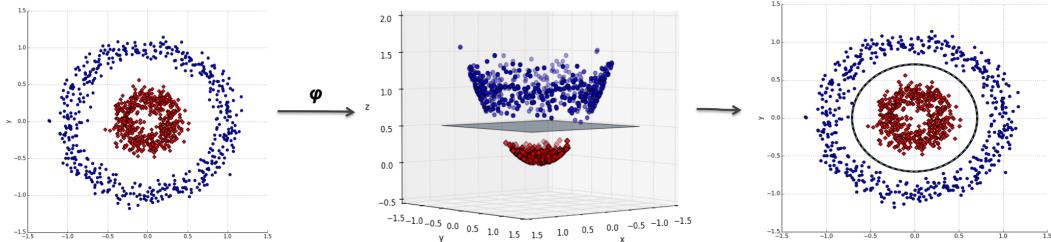


Figure 2.26. Transformation of the data in a feature space where the instances from the two classes may be linearly separable [92].

Chapter 3

Architecture

In this chapter we will discuss the architecture of our framework and the tools used for this thesis in Deception Detection. We will first explain the general architecture to perform deception detection, then we will introduce the OpenFace toolkit (3.2), and describe landmark detection (3.3), action unit recognition (3.4) and the deception detection process (3.5).

3.1 General Architecture

The general idea of this thesis is to classify whether a person is telling the truth or lying by analyzing the subject's facial movement, extracted using an RGB camera, during a video. To achieve this it is fundamental that the subject's face results visible during the videos.

We implemented this idea using a Support Vector Machine classifier.

The architecture utilized is shown in fig ?? and is composed this way:

- We gather a video of a person in a high stake situation.
- Detect the face of the person in the video using Constrained Local Neural Field.
- Detect the facial landmarks using Convolutional Experts Network.
- Extract the appearance features in the form of Histogram of Oriented Gradients and perform dimensionality reduction and feature normalization through Principal Component Analysis.
- Extract action unit presence and intensity from the videos using a SVM.
- Feed the extracted data to train an SVM for classification on a test set.

3.2 OpenFace

The OpenFace [93] toolkit is a tool for machine learning and computer vision researchers, created by Baltrušaitis et al., to support performing facial landmark detection, head pose estimation, action unit recognition, feature extraction and eye gaze estimation.

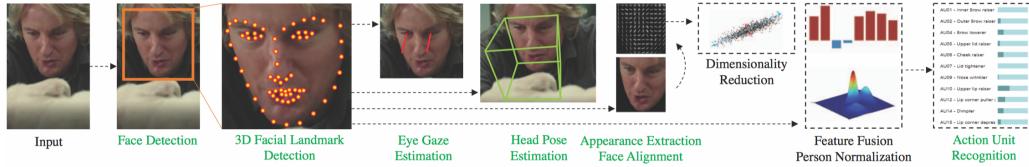


Figure 3.1. OpenFace Pipeline [93]

This tool is customizable in the way that all the SVM and neural networks models can be trained with different datasets, and the code is open source so it can be changed to fit the needs of the experimenters.

3.3 Landmark Detection

The first step in Action Unit identification is to detect the facial landmarks. To accomplish this a local detector called Convolutional Experts Network (CEN) (Fig. 3.2) is utilized.

CEN has the advantage of aggregating a neural architecture and patch experts which are local detectors and are used to evaluate the probability of a landmark being aligned at a particular pixel location. CEN can learn different patch experts and adapt to diverse appearance models without explicitly labeling attributes.

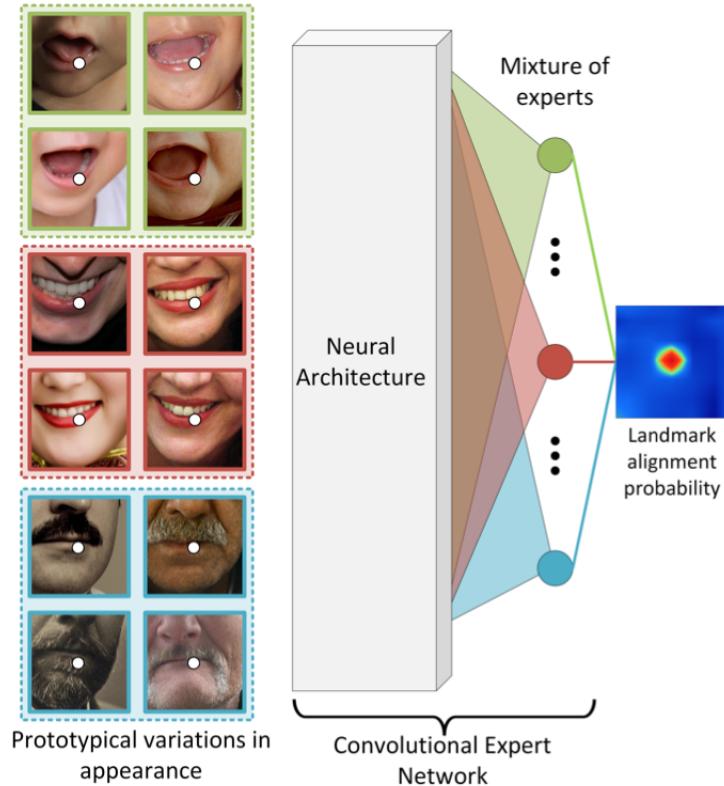


Figure 3.2. Facial landmarks naturally cluster around appearance prototypes (facial hair, expressions, make-up etc). To model such appearance variations a Convolutional Experts Network is used to merge the advantages of neural architectures and mixtures of patch experts, to model landmark alignment probability [109].

OpenFace uses a Convolutional Experts Constrained Local Model (CE-CLM) [109], which is a Constrained Local Model (CLM) that uses CEN as a local detector.

A Constrained Local Model is class of methods for locating sets of points (constrained by a statistical shape model) on a target image [25].

Generally the procedure is as follows (Fig. 3.3):

- Sample a region from the image around the current estimate, projecting it into a reference frame.
- For each point, generate a "response image" giving a cost for having the point at each pixel.
- Search for a combination of points which optimizes the total cost, by manipulating the shape model parameters.



Figure 3.3. Overview of CLM [25].

CLM are used to model the appearance of each facial landmark individually by using local detectors and a shape model for constrained optimization.

The CE-CLM is divided in two fundamental parts: response map computation using CEN, and shape parameter update.

In the first step, the landmarks alignment is computed separately from the other landmarks.

In the second phase, all landmarks are considered together and for misaligned landmarks and irregular shapes their position is penalized, using a Point Distribution Model (PDM).

3.3.1 CEN

In the initial step, the objective is to generate a response map to localize the individual landmarks. This is performed by assessing the landmark alignment probability at specific pixel locations.

CEN takes in input a Region of Interest (ROI) around the currently estimated position of a landmark, and outputs a response map that calculates the landmark alignment probability at each pixel location (Fig. 3.4).



Figure 3.4. Overview of the Convolutional Experts Network model. The output response map is a non-negative and non-linear combination of neurons in ME-layer using a sigmoid activation [109].

In order to do all this, the ROI is initially passed in a Contrast Normalizing Convolution layer to perform z-score normalization and to calculate the correlation between input and kernel. The output is then convolved in another layer of ReLU neurons (convolution is an operation on two functions to produce a third function that expresses how the shape of one is modified by the other).

The last neural layer before the response map is the Mixture of Expert Layer (ME-layer), and it can model the alignment probability using a combination of patch experts (local detectors) that are able to represent different landmarks appearance prototypes by outputting individual votes on alignment through a sigmoid function. The response maps from all the local detectors are then combined in the last layer, giving the final alignment probability.

3.3.2 Point Distribution Model

The Point Distribution Model is used to represent the mean geometry of a shape and some statistical modes of geometric variation inferred from a training set of shapes [104].

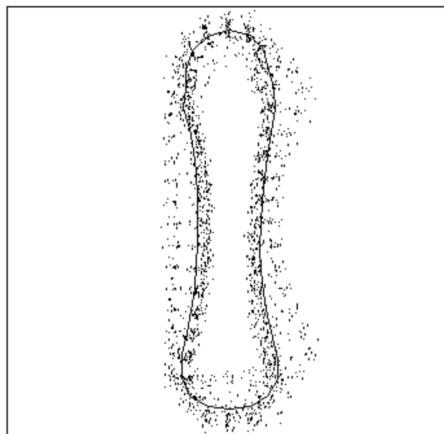


Figure 3.5. PDM of a metacarpal. Dots mark the possible position of landmarks, and the line denote the mean shape [47].

Point distribution models rely on landmark points. The general PDM works this way:

1. A set of training images are manually landmarked to sufficiently approximate the geometry of the original shapes.
2. This k landmarks are aligned in two dimensions resulting in $\mathbf{X} = (x_1, y_1, \dots, x_k, y_k)$
3. The shape outlines are reduced to sequences of k landmarks, so that any given training shape can be defined by the vector $\mathbf{X} \in \mathbb{R}^{2k}$.
4. The matrix of the top d eigenvectors is given as $\mathbf{P} \in \mathbb{R}^{2k \times d}$, and each eigenvector describes a principal mode of variation along the set.
5. A linear combination of the eigenvectors is used to define a new shape \mathbf{X}' , mathematically defined as:

$$\mathbf{X}' = \bar{\mathbf{X}} + \mathbf{P}\mathbf{b} \quad (3.1)$$

where $\bar{\mathbf{X}}$ is defined as the mean shape across all training images, and \mathbf{b} is a vector of scaling values for each principal component.

6. By modifying the variable \mathbf{b} an infinite number of shapes can be defined. \mathbf{b} shouldn't generally be modified more than $\pm 3\sigma$ [104].

With the OpenFace framework Point Distribution Model [82] we model the location of facial feature points in the image, using non-rigid shape and rigid global transformation parameters.

The application of PDM has two objectives:

- They are employed to control the landmark locations.
- They are used to regularize the shapes in the CE-CLM framework by using Non-Uniform Regularized Landmark Mean Shift (NU-RLMS) [13].

This has an effect in the final detected landmarks, where the irregular shapes are imposed a penalty.

The results can be seen in the figure 3.6 below for different images with different illuminations and angles.



Figure 3.6. Example of different facial landmarks detected in diverse conditions and viewing angles. [93].

3.4 Action Unit Detection

Action Unit (AU) detection plays a fundamental role in our work. We now describe the process used to detect the AUs through an SVM, starting with an overview of the datasets utilized for training the SVM in AUs detection.

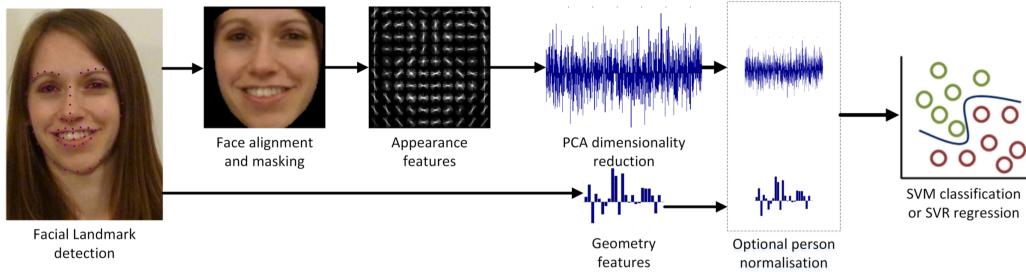


Figure 3.7. AU detection and intensity pipeline [12].

3.4.1 Training Datasets

There are four main dataset used for training the Action Unit detection system: DISFA [60], BP4D-Spontaneous [110], SEMAINE [62] and CK+ [57]. These four datasets consist of videos of people subject to emotion inducing tasks.

The **BP4D** database of spontaneous facial expressions includes videos of 41 participants (23 women and 18 men, 21 for training and 20 for validating). The age ranged from 18 to 29; 11 were Asian, 6 African-American, 4 Hispanic, and 20 Euro-American.

Emotion inducing techniques were used to elicit an emotional response. Frame-level ground-truth for facial actions was obtained by using the Facial Action Coding System annotations, performed by trained professionals.

Each participant in the database is associated with 8 tasks. For each task, there are both 3D and 2D videos and the metadata include annotations for 11 AUs for occurrence and 5 AUs for intensities.

DISFA (Denver Intensity of Spontaneous Facial Action) Database is a non posed facial expression database for automatic action unit detection. It contains videos of 27 participants (12 female, 15 males; 14 used for training and 13 for validation). It includes 4 minute-long videos of spontaneous facial expression, resulting in 130k frames annotated for 12 AUs (Fig. 3.8), comprehensive of AUs intensity on a 0 to 5 scale.

AU#	FACS Name
1	Inner Brow Raiser
2	Outer Brow Raiser
4	Brow Lowerer
5	Upper Lid Raiser
6	Cheek Raiser
9	Nose Wrinkler
12	Lip Corner Puller
15	Lip Corner Depressor
17	Chin Raiser & Mentalis
20	Lip Stretcher
25	Lip Part
26	Jaw Drop

Figure 3.8. AUs coded in the DISFA database [28].

The FACS annotated **SEMAINE** subset contains recordings of 31 subjects (15 for training and 16 for validation). It consists of one minute long recordings, leading to 93k frames labeled for 5 AU occurrences.

The **Cohn-Kanade** AU-Coded Facial Expression Database is utilized for research in facial image analysis.

CK+, includes both posed and non-posed (spontaneous) expressions and additional types of metadata. Each posed expression is completely FACS coded, and emotion labels are part of the metadata.

CK+ also provides baseline results for facial feature tracking, action unit and emotion recognition.

BP4D, SEMAINE and DISFA have three AUs in common (2, 12, and 17).

SEMAINE and DISFA share AUs 2, 12, 17, 25.

BP4D and DISFA share AUs 1, 2, 4, 6, 12, 15, 17.

This allows for cross-dataset training.

By training on these dataset is possible to detect the following AUs (3.9):

AU	Full name	Illustration
AU1	INNER BROW RAISER	
AU2	OUTER BROW RAISER	
AU4	BROW LOWERER	
AU5	UPPER LID RAISER	
AU6	CHEEK RAISER	
AU7	LID TIGHTENER	
AU9	NOSE WRINKLER	
AU10	UPPER LIP RAISER	
AU12	LIP CORNER PULLER	
AU14	DIMPLER	
AU15	LIP CORNER DEPRESSOR	
AU17	CHIN RAISER	
AU20	LIP STRETCHED	
AU23	LIP TIGHTENER	
AU25	LIPS PART	
AU26	JAW DROP	
AU28	LIP SUCK	
AU45	BLINK	

Figure 3.9. List of available AUs for prediction [93].

3.4.2 Feature Extraction

There are two types of features that are used: appearance (Histogram of Oriented Gradients) and geometry (shape parameters and landmark locations) ones. To extract those features it's required to track certain landmarks on the face, and then continue this process by performing face alignment.

Face Tracking

Face tracking is done by utilizing Constrained Local Neural Field (CLNF) facial landmark detector and tracker, backed up by a structural SVM for facial detection [13]. CLNF is a specific case of Constrained Local Model (CLM), that differs from

the original by utilizing more advanced local detectors and a different optimization function.

The Constrained Local Model can be described by the following parameters:
 $p = [s, \mathbf{R}, \mathbf{p}, \mathbf{t}]$.

- s is the scale factor.
- \mathbf{R} is the object rotation.
- \mathbf{t} is the 2D translation.
- \mathbf{p} describes the shape of a vector of non rigid variations.

These parameters can be modified to compute different versions of the model. The resulting point distribution model is:

$$x_i = s \cdot \mathbf{R}(\bar{x}_i + \phi_i \mathbf{p}) + \mathbf{t} \quad (3.2)$$

Where x_i is the location of the i th feature point in an image, \bar{x}_i is the mean of the i th element of the PDM, and ϕ_i is the i th eigenvector that describes the variation of the feature point.

In the Constrained Local Model we use the parameters from face detection to estimate the maximum a posteriori probability p of the face model.

Alignment and Masking

For the extracted face to be correctly analyzed, there needs to be a mapping to a common reference frame, and the changes resulting from scaling and in plane rotation needs to be removed.

In order to achieve this, a similarity transform from the currently detected landmarks to a representation of frontal landmarks from a neutral expression (projection of mean shape from a 3D PDM) is utilized.

The similarity transform is done with Procrustes superimposition that minimizes the mean square error between aligned pixels [13].

The result is a 112×112 pixel image of the face with 45 pixel interpupillary distance (Fig 3.10).

To reduce the weight of significant facial expressions (mouth opening, brow raises etc.) on the similarity transform, only the most stable facial landmarks must be used.

In order to determine those points, the most stable CLNF detected landmarks on the CK+ dataset [57] are examined.



Figure 3.10. Stable points used for alignment to a common reference frame, followed by masking [12].

Masking, performed using a convex hull surrounding the feature points, is done to remove informations from the image that do not regard the face.

Appearance Features

After the face is aligned to a 112×112 image it's time to extract the appearance features.

To get the appearance features, 12×12 block of 31 dimensional Histogram of Oriented Gradients (HOG) are extracted, giving a 4464 dimensional vector characterizing the face. The implementation for HoG comes from dlib [50].

Once the feature vector is obtained, the next step is to reduce it's dimensionality. In order to do that Principal Component Analysis (PCA) is applied.

To generalize the dimensionality reduction, the training was performed on the FERA 2015 [94], CK+ [57], DISFA [60], AVEC 2011 [83] and FERA 2011 [95] datasets.

By using PCA while sub-sampling from peak and neutral expressions and keeping 95% of explained variability, the feature vector becomes of 1379 dimensions.

Geometry Features

The geometry features are obtained trough the CE-CLM models, and consist of non-rigid shape parameters and landmarks locations (p and $\phi_i p$ in equation 3.2). this results in a 227 dimensional vector representing geometry features.

Summed to the appearance features, this leads to a 1606 dimensional vector that defines the appearance of the face.

Neutral Expression Extraction

To extract some of the facial expression is very important to have a neutral starting expression. There are personal differences on how we appear while in a neutral resting state, such as people looking naturally more cheerful or sad then others

[66]. To address this issue there needs to be a person specific calibration, done by adjusting for neutral expressions [13].

Neutral expression adjustments are done by computing the median value of the face descriptors in a video, leading to a neutral expression descriptor. This works assuming that in a video most of the frames will contain a neutral expression, and this should hold true especially for real life situations where most of the time the interactions are performed with a neutral expression [8].

Once the neutral descriptor is computed, it is subtracted from the feature descriptor, giving normalized features.

Classification and Regression

The next step is to extract Action Units from each frame of the video. Both intensity and presence are extracted.

Action Units detection is performed through Support Vector Machines (SVM), and Action Unit intensity using Support Vector Regression (SVR), both using the liblinear implementation [65]. In both cases a linear kernel is utilized, as more complex kernels had no effect on performances and were quite slower to train.

Since the AU occurrences are not balanced by nature (for example AU28, blinking, is not very frequent if we analyze each frame of a video), it's highly important to balance the training data. To do this it was needed to undersample the negative AU samples from training data, thus getting to an equal number of positive and negative samples.

3.5 Deception Detection

In this step we feed the extracted action units from the videos to an SVM for classification. We use an SVM because after a lot of testing with different techniques that was the one that performed the best compared to others.

3.5.1 Data pre-processing

Before feeding the videos to the SVM for Action Units extraction we had to make sure the videos were in good enough condition to be processed. That is not given since all this videos come from real life trials and the recording conditions are not optimal.

Being in an acceptable condition means that the face should be clearly visible, and the video has to have the specific person visible and talking in at least *most* of the video, and no other people on it, so to avoid training on different subjects of which we are not sure if they are lying or telling the truth.

As a consequence many videos had to be cut, trimmed or discarded to reach an acceptable level.

After extracting Action Units intensity and presence from our dataset of videos (4.2) we put the data in a *csv* format to be able to process them and pass them to the SVM.

The data are divided for both AUs intensity, ranging from 0 to 5, and AU presence, having the boolean value of either 0 or 1 if the AU was missing or present in the specific frame of the video.

AU12_r	AU14_r	AU15_r	AU17_r	AU20_r	AU23_r	AU25_r	AU26_r	AU45_r	AU01_c	AU02_c	AU04_c	AU05_c	AU06_c	AU07_c	AU09_c	AU10_c	AU12_c	AU14_c
1.33	2.2	0.22	0.12	0	0	0.53	0	0.83	0	0	1	0	0	0	0	1	1	1
1.31	2.11	0.22	0.04	0.02	0	0.47	0	0.73	0	0	0	0	0	0	0	1	1	1
1.3	2.07	0.28	0	0.02	0	0.35	0	0.65	0	0	0	0	0	0	0	1	1	1
1.28	2.11	0.31	0	0	0	0.26	0	0.58	0	0	0	0	0	0	0	1	1	1
1.27	2.11	0.42	0.03	0	0	0.21	0	0.56	0	0	0	0	0	0	0	1	1	1
1.26	2.12	0.37	0.04	0	0	0.28	0	0.5	0	0	0	0	0	0	0	1	1	1
1.31	2.08	0.45	0.13	0	0.01	0.36	0	0.55	0	0	0	0	0	0	0	1	1	1
1.44	2.05	0.47	0.26	0	0.02	0.26	0	0.61	0	0	0	0	0	0	0	1	1	1
1.55	2.04	0.51	0.35	0	0.02	0.18	0	0.68	0	0	0	0	0	0	0	1	1	1
1.56	2.01	0.43	0.38	0	0.01	0.12	0	0.69	0	0	0	0	0	0	0	1	1	1
1.53	2.06	0.3	0.36	0	0	0.17	0	0.67	0	0	0	0	0	0	0	1	1	1
1.52	2.08	0.27	0.38	0	0	0.18	0	0.71	0	0	0	0	0	0	0	1	1	1
1.55	2.14	0.19	0.38	0	0	0.27	0	0.72	0	0	0	0	0	0	0	1	1	1
1.59	2.1	0.23	0.34	0	0	0.33	0	0.71	0	0	0	0	0	0	0	1	1	1
1.6	2.11	0.32	0.36	0	0.03	0.37	0	0.67	0	0	0	0	0	0	0	1	1	1

Figure 3.11. Example of input data in csv format.

The next step is to remove all the rows that have both no intensity and no presence for AUs, since those rows are most likely errors from the extraction, deriving from the poor quality of the videos, or from the face not being recognized in some frames.

We organized the data in three sets, one concerning only AUs intensity, one with only AUs presence, and another one where we merged the previous ones together, in order to better analyze the data and understand its structure.

3.5.2 SVM

After pre-processing the data we use the train set to train a sigmoid SVM, and then we utilize the trained model to evaluate on the test set and compare the results with the labels present on the test set.

Generally, each of n frames is converted in a vector of length 36, where 17 variables represent the intensity and 18 variables represent the presence of Action Units. The last variable is the class label.

So a given training input \mathbf{AU} for our SVM, \mathbf{AU} is a $36 \times n$ matrix composed of AUI_i and AUP_i , where AUI_i represents the intensity for action units in row i and AUP_i is the action units presence for row i .

These observations belong to two classes, truthful and deceptive, represented by a label $\mathbf{Y} = \{y_1 \dots y_n\}$, where $y = \{0, 1\}$ depending on whether the frame belongs to a deceptive or truthful observation.

$$\mathbf{AU} = (AUI_i + AUP_i, y_i) \quad (3.3)$$

An input vector au^* for the test set, this time of length 35, is formed the same way but without having a class label.

$$\mathbf{au}^* = (AUI_1^* + AUP_1^*, \dots, AUI_{35}^* + AUP_{35}^*)^T \quad (3.4)$$

The goal of our SVM classifier is to correctly classify the test observation based on the training done by using the data in the train set.

Suppose we are using a *linear* SVM classifier. We would need to figure out whether a point lies on one side or the other of an hyperplane, and we would look at this kind of equations to classify a test observation au^* based on the sign of $f(au^*)$:

$$f(au^*) = \beta_0 + \beta_1 au_1^* + \beta_2 au_2^* + \dots + \beta_{35} au_{35}^* \quad (3.5)$$

if $f(au^*)$ is positive, then we would assign the test observation to class 1, and if it is negative, we would assign it to class 0.

Of course there are two things we need to consider: - we need to keep the margin M of the hyperplane in consideration - the relationship between features is not linear so we could need to enlarge the feature space

For the Support Vector Machine we enlarge the feature space using kernel functions [2.4.6](#)

Chapter 4

Experiments

This describes the stack used for this thesis and the experiments done to achieve our results: I will start by reviewing the configurations and libraries utilized (par. 4.1), the database used to acquire training and test set (par. 4.2), all the analysis done on the data (par. 4.3), and finally we will discuss the results of this work (par. 4.4).

4.1 Building Blocks

To realize this work we used the following tools and libraries:

- OpenFace: the library to train the models for face detection, landmark detection, feature extraction and action unit recognition.
- R + R-Studio: the environment where I implemented the code, with some important libraries, such as:
 - e1071: library for SVM classification (implementation of libsvm).
 - Random Forest: library to perform Random Forest and calculate variable importance.
 - Corrplot: library for calculating and visualizing correlations.

The experiments were conducted on two systems, Windows and Mac OS, with the following specs:

The windows machine has 8GB of RAM, i5 processor, and 4096MB ATI AMD Radeon R9 GPU.

The Mac has 4gb RAM, i5 processor, and an intel integrated GPU.

SVM training took about one hour on both machines.

4.2 Real Life Trial Database

This section is about the database we used to perform our experiments: it comes from the work done for the paper "Deception Detection using Real-life Trial Data" [71].



Figure 4.1. Examples of images from the dataset videos. [71].

The dataset is gathered from real-life trial videos available on YouTube and other public websites. The dataset also contains statements made by exonerees after exoneration, and some statements from defendants during crime-related TV episodes.

The first step to collecting the dataset was to identify public multimedia sources where the recordings of the trials were available, and deceptive and truthful behavior could be observed and verified.

The videos are of trial recordings where the defendant or witness in the video can be clearly identified, the face is visible enough during most of the clip duration, and the visual quality should be good enough to accurately see the facial expressions (Fig. 4.1).

There are three outcomes for the trials that were considered to label the videos as deceptive or truthful: guilty, non-guilty, and exoneration.

For the guilty verdicts the deceptive clips are taken from the defendant in the trial, while the truthful clips are gathered from the witnesses. There are also instances where the deceptive videos are of suspects denying a committed crime, and truthful ones are from the same person answering questions that were verified by the police as truthful.

In regards to the witnesses, if the testimony is verified by a police officer they are labeled as true.

Testimonies that help the guilty party are labeled as false. Exoneration (reversal of the sentence) testimonies are regarded as truthful.

The original dataset consists of 121 videos, 61 of which are deceptive and 60 truthful.

The average length of the videos is 28.0 seconds. The average video length for deceptive videos is 27.7 seconds, while the one for truthful videos is 28.3 seconds. The data consists of 58 total subject, 22 females and 36 males, with ages between 16 and 60 years.

This dataset was annotated following the MUMIN coding scheme for hand movement and facial displays. While annotating, the annotators were able to chose only one label per gesture, for every clip.

Fig. 4.2 shows the frequency counts associated with the gestures considered during the annotation.

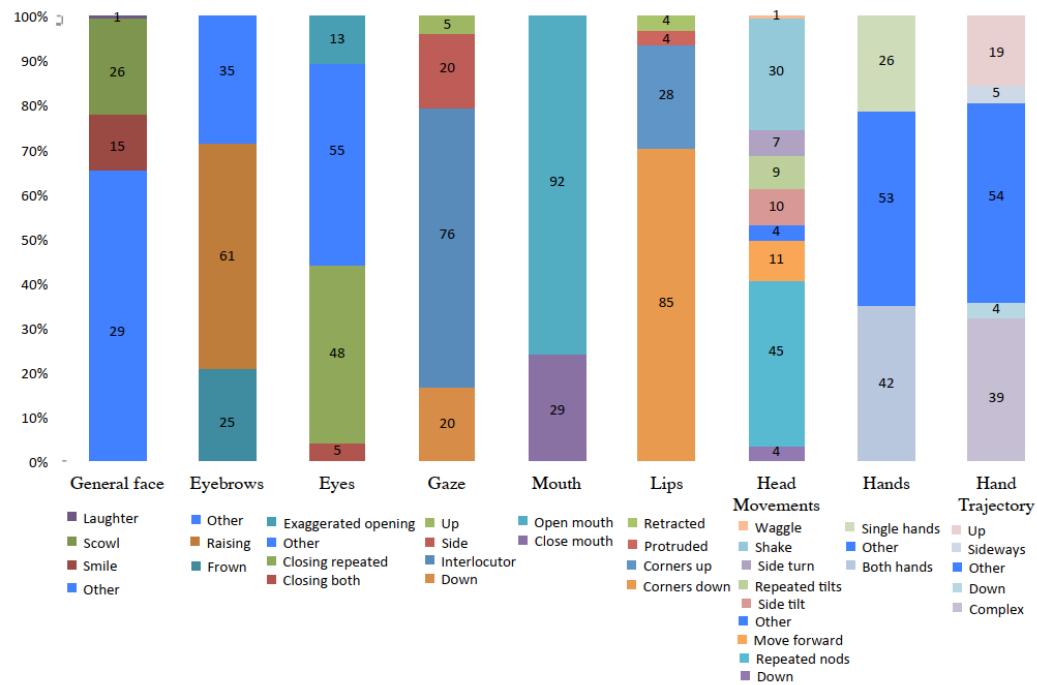


Figure 4.2. Frequency of gestures annotated in the videos from the database. [71].

We modified this database in the following way:

- We cut or removed videos where the where the face was covered, hidden or in any way very difficult to see or recognize.
- We also cut parts of video that had multiple subjects in the scene because Action Unit extraction works on one person only at the moment.
- Since it is necessary to see the person to perform AU extraction, we also avoided some videos in which the subject was not visible while talking, but the interlocutor was shown.

With deceptive videos going from #1 to #61 and truthful videos going from #1 to #60 we cut and removed the following videos from the original dataset, since they showed signs of the problems just described:

- **Deceptive:**

- CUT: Video numbers 46, 47, 48, 49, 52, 54, 56.
 - REMOVED: Video numbers 50, 53, 55.

- **Truthful:**

- CUT: Video number 7, 28, 43.
 - REMOVED: Video numbers 12, 31.

We also performed a division of subjects in the training and test set "by hand", splitting data to around 75% for the training set and 25% for the test set, to train the classifiers in a way that it wouldn't adapt to specific people, a problem that could occur since there are many videos with the same subject. In fact, the same person never appears both in the training and in the test set as an effort to completely avoid this problem.

The training set consists of 72451 observations of 36 variables, and the test set is composed of 14198 observations of 36 variables.

4.3 Data Analysis

In the following sections we describe the methods used to analyze the data and the machine learning techniques to classify them.

We start with comparing the occurrences for deceptive and truthful videos on the training set (par. 4.3.1), then look at the correlation between variables (par. 4.3.2), then we explain how we used GLM (par. 4.3.3), Random Forest (par. 4.3.4) and Support Vector Machines (par. 4.3.5)

4.3.1 Data Comparison

The first thing I did, after cleaning the data, was to compare the extracted deceptive (Fig 4.3) and truthful (Fig. 4.4) AUs occurrences (only for presence, not intensity) from the training set, which is shown in Fig. 4.5.

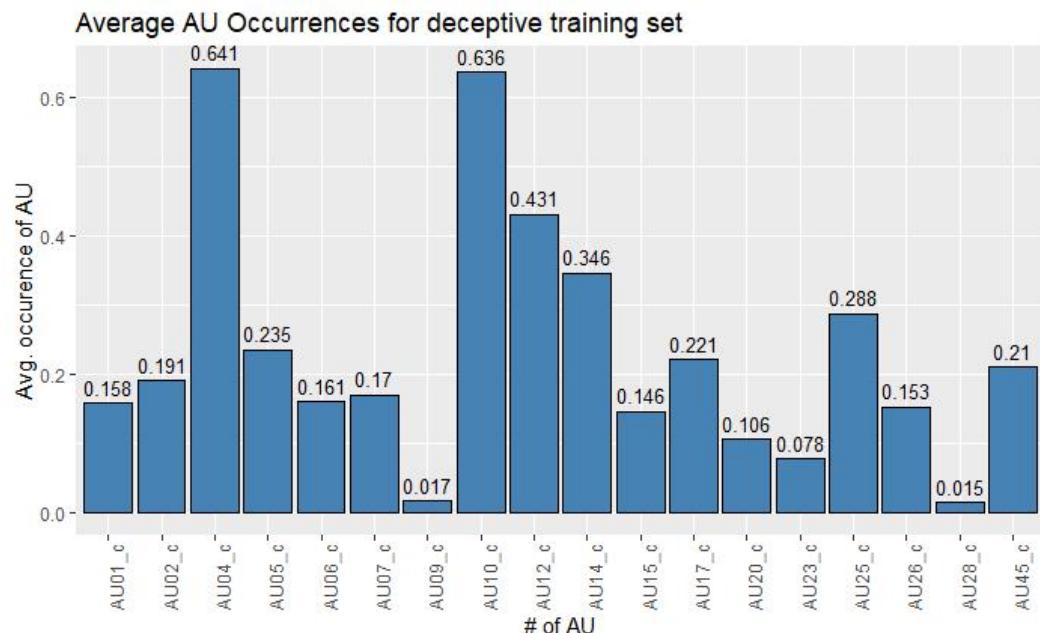


Figure 4.3. Average AU Occurrences for the deceptive training set.

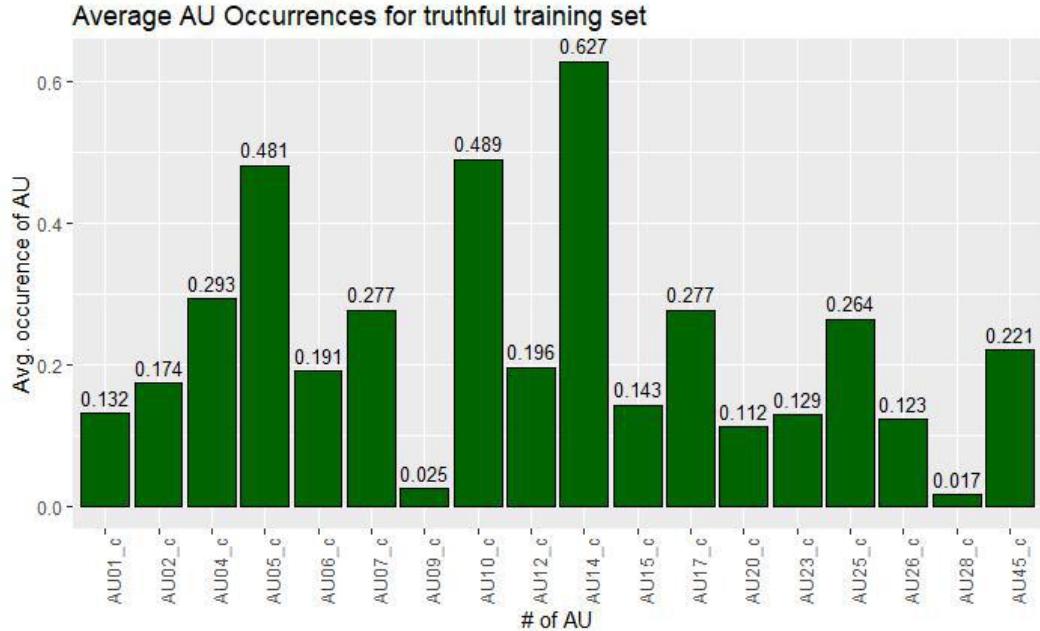


Figure 4.4. Average AU Occurrences for the truthful training set.

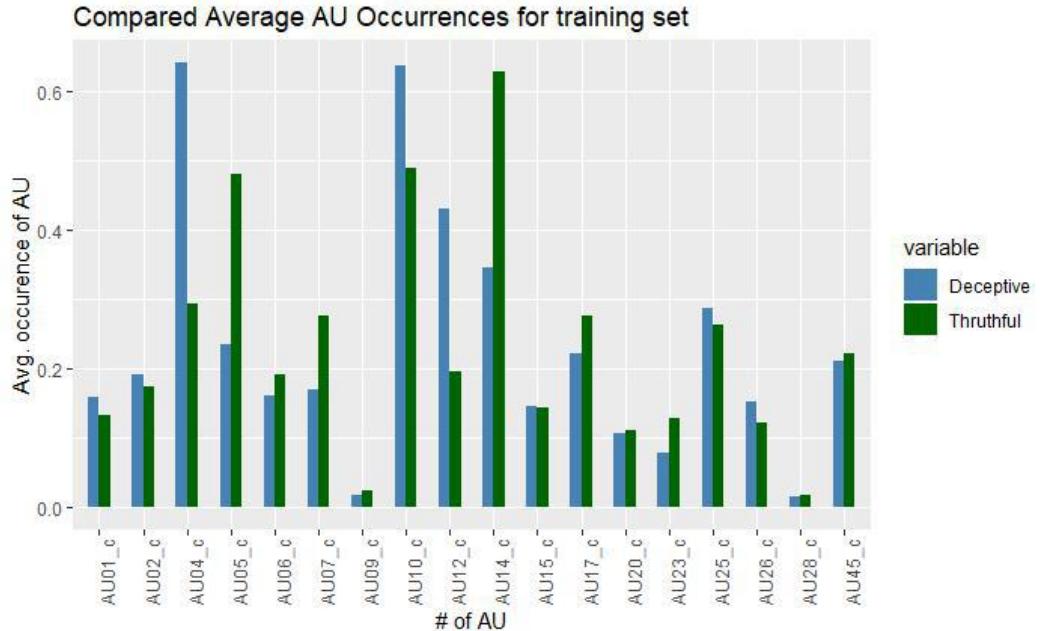


Figure 4.5. Comparison of AU Occurrences for the training set.

This comparison shows interesting differences between truthful and deceptive occurrences for AU04 (Brow Lowerer), AU05 (Upper Lid Raiser), AU10 (Upper Lip Raiser), AU12 (Lip Corner Puller) and AU14 (Dimpler).

4.3.2 Correlation

Variables correlation is shown in figure 4.6. When the AUs label ends with "_r" it indicate intensity and "_c" means presence.

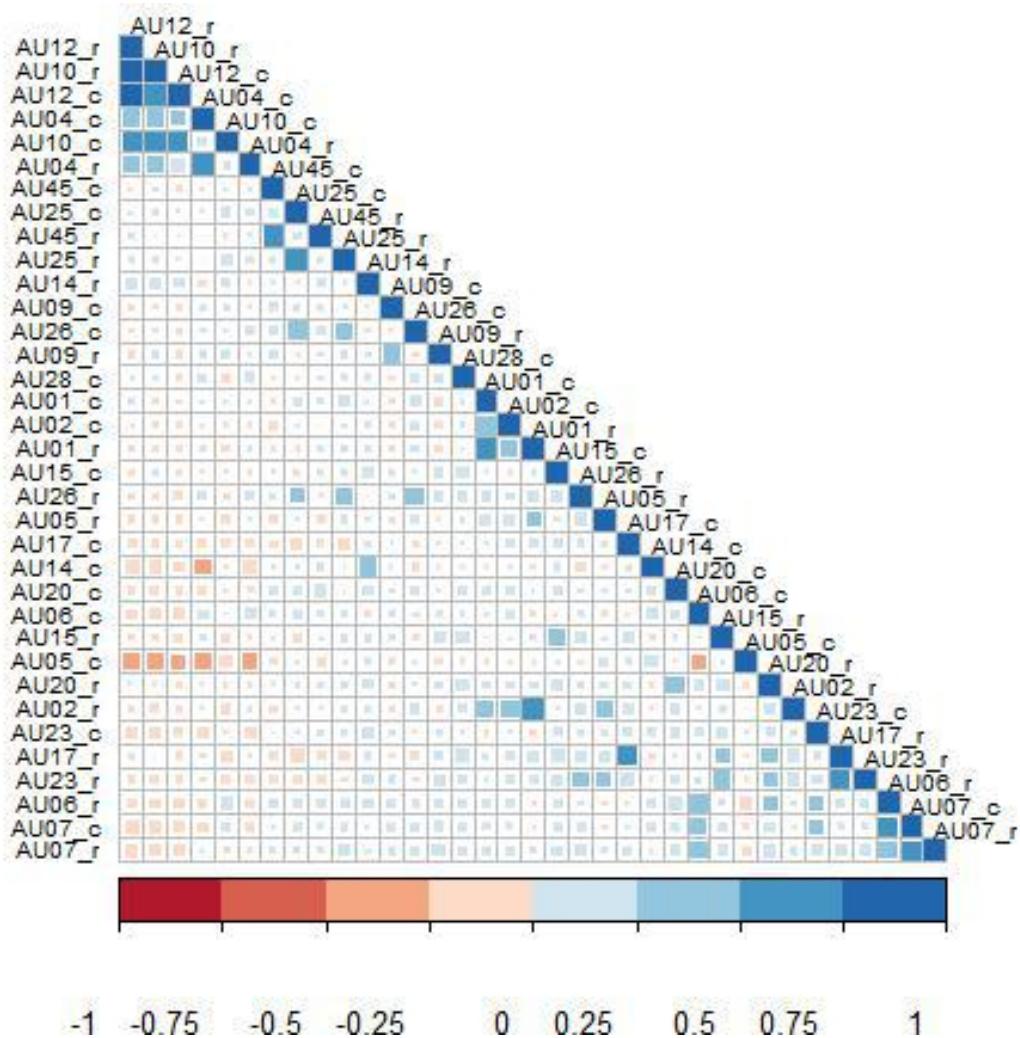


Figure 4.6. Correlation between variables for both presence and intensity in the training set.

which shows some obvious correlations between presence and intensity of the same AU, and significant correlation between AU10_r, AU12_r, AU12_c, AU07_r, AU17_r, AU06_r, AU04_c, AU02_r, AU01_r, AU25_r, AU45_c.

4.3.3 GLM

Another quick experiment I made was with the Generalized Linear Model in R, which translates into Logistic Regression (Par. 2.2), at least to have an idea of the resulting p-values, as seen in fig. 4.7:

	Pr(> z)
(Intercept)	2.58e-15 ***
AU01_r	0.60606
AU02_r	1.14e-06 ***
AU04_r	1.21e-12 ***
AU05_r	7.80e-12 ***
AU06_r	3.14e-09 ***
AU07_r	< 2e-16 ***
AU09_r	< 2e-16 ***
AU10_r	< 2e-16 ***
AU12_r	< 2e-16 ***
AU14_r	< 2e-16 ***
AU15_r	< 2e-16 ***
AU17_r	2.46e-13 ***
AU20_r	< 2e-16 ***
AU23_r	0.83786
AU25_r	< 2e-16 ***
AU26_r	3.88e-10 ***
AU45_r	< 2e-16 ***
AU01_c	1.59e-07 ***
AU02_c	< 2e-16 ***
AU04_c	< 2e-16 ***
AU05_c	< 2e-16 ***
AU06_c	6.30e-08 ***
AU07_c	0.00136 **
AU09_c	2.77e-15 ***
AU10_c	< 2e-16 ***
AU12_c	< 2e-16 ***
AU14_c	< 2e-16 ***
AU15_c	< 2e-16 ***
AU17_c	0.76875
AU20_c	< 2e-16 ***
AU23_c	< 2e-16 ***
AU25_c	< 2e-16 ***
AU26_c	< 2e-16 ***
AU28_c	8.89e-08 ***
AU45_c	< 2e-16 ***

Figure 4.7. p-values from GLM for each AU (intensity and presence).

As we can see AU01_r (Inner Brow Raiser), AU23_r (Lip Tightener), and AU17_c (Chin Raiser) are not significant and that's an indication that we could remove them from the features used to make the classification.

The results for classifying data with GLM is 54.86% accuracy on the test set, which is a little better than chance but not great.

4.3.4 Random Forest

I've decided to use Random Forest (par. 2.3) to estimate variable importance, shown in fig. 4.8.

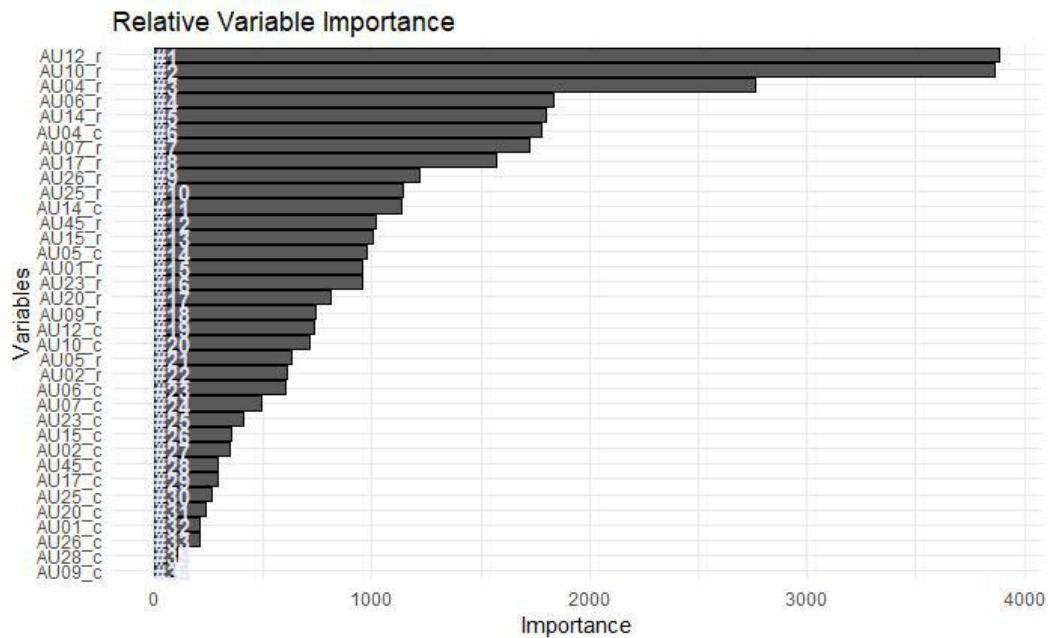


Figure 4.8. Variable Importance extracted from classifying data with Random Forest.

The number of grown trees are 200 as going further did not increase the prediction accuracy at all.

Fig. 4.8 shows that the intensity of AUs 04, 06, 07, 10, 12 and 14 are the most significant, while the sole presence of an AU is not particularly significant, beside for AUs 04, 05 and 14.

As we can see AUs that have low p-values in GLM have also low importance in Random Forest.

Classifying with Random Forest gave me 58,42% accuracy on the test set, which is better than GLM but it is still not good enough to be significant.

4.3.5 SVM

To approach classification with an SVM we tried both linear, radial and sigmoid kernel.

The cost c and γ parameter were found using a tuning function that performs a grid search over the combinations of cost and gamma for all the kinds of kernels we used, and then by empirically changing the values to see if better results could be achieved.

In the tuning function, γ varied from 10^{-5} to 10^{-1} and c from 10^{-3} to 10^2 , growing by 10^1 every iteration.

After the grid search I did some experimental changes and found empirically that the best parameters for the linear kernel were $C = 1$ and $\gamma = 0.04$

The results for classifying with a radial kernel are 67% accuracy on the test set (to be precise 66.99394), which is a good starting result, while the linear kernel performed worse at around 61% accuracy.

The best scoring kernel was actually the sigmoid kernel, giving us 72,64206% accuracy when predicting the test set, by setting $\gamma = 0.001$ and $C = 3$, obtaining the values of the parameters via both the tuning function and trial and error.

4.4 Results

Using Support Vector Machine we were able to classify deceptive and truthful videos with an accuracy of 72,64%.

The table below shows a comparison of the results obtained using different machine learning algorithms, and different kernels for SVM.

Algorithm	Accuracy
Logistic Regression	54,86
Random Forest	58,42
SVM - Linear	61,56
SVM - Radial	66,99
SVM - Sigmoid	72,64

Table 4.1. Comparison of the results obtained with different approaches for this dataset

Chapter 5

Conclusions

5.1 Final Considerations

In this thesis we presented a way to discriminate between lies and truths, based on the extraction of AUs presence and intensity from videos, through an SVM classifier.

The applications for a system that can accurately detect detection are numerous and space in a lot of fields: from airport security to job interviews.

5.2 Future Work

Some possible developments to be applied to this thesis are:

- One of the difficulties of this kind of studies is the lack of a big dataset for high stake deception. It would be beneficial to gather such dataset to improve the results, bettering also the quality of the videos.
- Extend AUs extraction to more than one person in the scene. It would be interesting to consider direct interaction with questions and responses.
- Using temporal information for each video to capture the beginning and the end of the interaction. The lie does not necessarily appear in *all* the video, but probably just at some point, even though the majority of the statements are either deceptive or truthful. We were thinking about LSTM to observe the temporality.
- More modules to consider other than the face. For example we could consider the body, and analyze the voice and speech patterns to extract information from the subject.

Bibliography

- [1] URL: <https://www.slideshare.net/mlvlc/14-ensembles-of-decision-trees> (p. 36).
- [2] URL: <https://slideplayer.com/slide/10391576/> (p. 37).
- [3] Sean A. Spence. “Playing Devil’s advocate: The case against fMRI lie detection”. In: 13 (Feb. 2008), pp. 11–25 (p. 10).
- [4] M. Abouelenien et al. “Detecting Deceptive Behavior via Integration of Discriminative Features From Multiple Modalities”. In: *IEEE Transactions on Information Forensics and Security* 12.5 (May 2017) (p. 12).
- [5] Mohamed Abouelenien, Rada Mihalcea, and Mihai Burzo. “Analyzing Thermal and Visual Clues of Deception for a Non-Contact Deception Detection Approach”. In: *Proceedings of the 9th ACM International Conference on PErvasive Technologies Related to Assistive Environments*. 2016. URL: <http://doi.acm.org/10.1145/2910674.2910682> (p. 12).
- [6] Mohamed Abouelenien, Rada Mihalcea, and Mihai Burzo. “Trimodal Analysis of Deceptive Behavior”. In: *Proceedings of the 2015 ACM on Workshop on Multimodal Deception Detection*. 2015. URL: <http://doi.acm.org/10.1145/2823465.2823470> (p. 11).
- [7] Mohamed Abouelenien et al. “Deception Detection Using a Multimodal Approach”. In: *Proceedings of the 16th International Conference on Multimodal Interaction*. 2014. URL: <http://doi.acm.org/10.1145/2663204.2663229> (p. 11).
- [8] S. Afzal and P. Robinson. “Natural affect data — Collection and annotation in a learning context”. In: *2009 3rd International Conference on Affective Computing and Intelligent Interaction and Workshops*. Sept. 2009. DOI: [10.1109/ACII.2009.5349537](https://doi.org/10.1109/ACII.2009.5349537) (p. 61).
- [9] Baker Alysha, ten Brinke Leanne, and Porter Stephen. “Will get fooled again: Emotionally intelligent people are easily duped by high-stakes deceivers”. In: *Legal and Criminological Psychology* 18.2 (), pp. 300–313. DOI: [10.1111/j.2044-8333.2012.02054.x](https://doi.org/10.1111/j.2044-8333.2012.02054.x). URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.2044-8333.2012.02054.x> (p. 3).
- [10] Tarek Amr. *Predict the Future with Regression Analysis*. URL: <https://medium.com/@gr33ndata/learn-regressions-analysis-23b789bf2c36> (p. 26).

- [11] A. Arasteh, M. H. Moradi, and A. Janghorbani. “A Novel Method Based on Empirical Mode Decomposition for P300-Based Detection of Deception”. In: *IEEE Transactions on Information Forensics and Security* 11.11 (Nov. 2016) (p. 9).
- [12] T. Baltrušaitis, M. Mahmoud, and P. Robinson. “Cross-dataset learning and person-specific normalisation for automatic Action Unit detection”. In: *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*. Vol. 06. May 2015, pp. 1–6 (pp. 16, 56, 60).
- [13] T. Baltrušaitis, P. Robinson, and L. Morency. “Constrained Local Neural Fields for Robust Facial Landmark Detection in the Wild”. In: *2013 IEEE International Conference on Computer Vision Workshops*. Dec. 2013 (pp. 54, 58, 59, 61).
- [14] T. Baltrušaitis, P. Robinson, and L. P. Morency. “OpenFace: An open source facial behavior analysis toolkit”. In: *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*. Mar. 2016, pp. 1–10 (p. 16).
- [15] Anup Bhande. *What is underfitting and overfitting in machine learning and how to deal with it*. URL: <https://medium.com/greyatom/what-is-underfitting-and-overfitting-in-machine-learning-and-how-to-deal-with-it-6803a989c76> (p. 23).
- [16] N. Bhaskaran et al. “Lie to Me: Deceit detection via online behavioral learning”. In: *Face and Gesture 2011*. Mar. 2011 (p. 8).
- [17] *Bias-Variance Tradeoff*. URL: <https://elitedatascience.com/bias-variance-tradeoff> (p. 22).
- [18] Marilyn G. Boltz, Rebecca L. Dyer, and Anna R. Miller. “Jo Are You Lying to Me? Temporal Cues for Deception”. In: *Journal of Language and Social Psychology* 29.4 (2010), pp. 458–466. DOI: [10.1177/0261927X10385976](https://doi.org/10.1177/0261927X10385976). URL: <https://doi.org/10.1177/0261927X10385976> (p. 5).
- [19] D. R. Burnham K. P.; Anderson. *Model Selection and Multimodel Inference (2nd ed.)*, Springer-Verlag. 2002 (p. 22).
- [20] Jr. Charles F. Bond and Bella M. DePaulo. “Accuracy of Deception Judgments”. In: *Personality and Social Psychology Review* 10.3 (2006), pp. 214–234. DOI: [10.1207/s15327957pspr1003_2](https://doi.org/10.1207/s15327957pspr1003_2). URL: https://doi.org/10.1207/s15327957pspr1003_2 (p. 3).
- [21] W. S. Chu, F. D. L. Torre, and J. F. Cohn. “Selective Transfer Machine for Personalized Facial Action Unit Detection”. In: *2013 IEEE Conference on Computer Vision and Pattern Recognition*. June 2013 (p. 17).
- [22] W. S. Chu, F. D. l. Torre, and J. F. Cohn. “Selective Transfer Machine for Personalized Facial Expression Analysis”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39.3 (Mar. 2017). DOI: [10.1109/TPAMI.2016.2547397](https://doi.org/10.1109/TPAMI.2016.2547397) (p. 17).

- [23] W. S. Chu, F. De la Torre, and J. F. Cohn. “Learning Spatial and Temporal Cues for Multi-Label Facial Action Unit Detection”. In: *2017 12th IEEE International Conference on Automatic Face Gesture Recognition (FG 2017)*. May 2017 (p. 16).
- [24] Wikipedia contributors. *Machine learning — Wikipedia, The Free Encyclopedia*. [Online; accessed 16-July-2018]. 2018. URL: https://en.wikipedia.org/w/index.php?title=Machine_learning&oldid=849817385 (p. 19).
- [25] Timothy Cootes. *Constrained Local Model*. URL: <https://personalpages.manchester.ac.uk/staff/timothy.f.cootes/Models/clm.html> (pp. 51, 52).
- [26] Langleben Daniel D. “Detection of deception with fMRI: Are we there yet?” In: *Legal and Criminological Psychology* 13.1 (), pp. 1–9. DOI: [10.1348/135532507X251641](https://doi.org/10.1348/135532507X251641). URL: <https://onlinelibrary.wiley.com/doi/abs/10.1348/135532507X251641> (p. 10).
- [27] Janey Davies. URL: <https://www.learning-mind.com/eye-movements-when-lying/> (p. 7).
- [28] *Denver Intensity of Spontaneous Facial Action Database*. URL: <http://www.engr.du.edu/mmahoos/DISFAContent.htm> (p. 57).
- [29] Shichuan Du, Yong Tao, and Aleix M. Martinez. “Compound facial expressions of emotion”. In: *Proceedings of the National Academy of Sciences of the United States of America* (2014) (p. 16).
- [30] P. Ekman and W. Friesen. “Facial Action Coding System: A Technique for the Measurement of Facial Movement”. In: *Consulting Psychologists Press* (1978) (p. 15).
- [31] P. Ekman and W. V. Friesen. “Nonverbal leakage and clues to deception”. In: *Psychiatry* 32.1 (1969), pp. 88–106 (p. 14).
- [32] Paul Ekman. *Emotions Revealed: Recognizing Faces and Feelings to Improve Communication and Emotional Life*. 2007 (p. 13).
- [33] Skrondal A. Everitt B.S. “Cambridge Dictionary of Statistics”. In: *Cambridge University Press* (2010) (p. 22).
- [34] Wittmann F. *Compare SVM Kernels*. URL: <https://gist.github.com/WittmannF/60680723ed8dd0cb993051a7448f7805> (p. 47).
- [35] Martha J. Farah et al. “Functional MRI-based lie detection: scientific and societal challenges”. In: *Nature Reviews Neuroscience* 15 (Jan. 2014). URL: <http://dx.doi.org/10.1038/nrn3665> (p. 10).
- [36] Kyosuke Fukuda. “Eye blinks: new indices for the detection of deception”. In: *International Journal of Psychophysiology* 40.3 (2001), pp. 239–245 (p. 7).
- [37] Trevor Hastie Gareth James Daniela Witten and Robert Tibshirani. *An Introduction to Statistical Learning*. URL: <http://www-bcf.usc.edu/~gareth/ISL/index.html> (pp. 28, 35–42, 44, 46).
- [38] S. George et al. “Eye blink count and eye blink duration analysis for deception detection”. In: *2017 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*. Sept. 2017 (p. 7).

- [39] E. Haggard and K. Isaacs. “Micromomentary facial expressions as indicators of ego mechanisms in psychotherapy”. In: *Methods of research in psychotherapy*. New York: Appleton-Century-Crofts (1966), pp. 154–165 (p. 14).
- [40] Ben Hammer. *Introducing Kaggle Scripts*. URL: <https://www.kaggle.com/general/13285> (p. 39).
- [41] Jeffrey Hancock. “Digital deception: When, where and how people lie online”. In: (Jan. 2012), pp. 287–301 (p. 3).
- [42] L. Hao et al. “Facial Action Unit Recognition Augmented by Their Dependencies”. In: *2018 13th IEEE International Conference on Automatic Face Gesture Recognition (FG 2018)*. May 2018 (p. 16).
- [43] Maria Hartwig et al. “Detecting deception in suspects: verbal cues as a function of interview strategy”. In: *Psychology, Crime & Law* 17.7 (2011), pp. 643–656. DOI: [10.1080/10683160903446982](https://doi.org/10.1080/10683160903446982). URL: <https://doi.org/10.1080/10683160903446982> (p. 3).
- [44] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning* (2nd ed.) Springer. 2008. URL: <http://www-stat.stanford.edu/~tibs/ElemStatLearn/> (pp. 33, 35).
- [45] Carl-Herman Hjortsjö. *Man's face and mimic language*. 1969 (p. 15).
- [46] Tin Kam Ho. “Random Decision Forests”. In: *Proceedings of the 3rd International Conference on Document Analysis and Recognition, Montreal, QC, 14–16 August 1995*. 1995 (p. 33).
- [47] *Image understanding: Point distribution models*. URL: <http://user.engineering.uiowa.edu/~dip/lecture/Understanding3.html> (p. 53).
- [48] *k-means data clustering*. URL: <https://towardsdatascience.com/k-means-data-clustering-bce3335d2203> (p. 24).
- [49] Michal Kawulok et al. “In Search of Truth: Analysis of Smile Intensity Dynamics to Detect Deception”. In: *Advances in Artificial Intelligence - IBERAMIA 2016*. 2016 (p. 14).
- [50] Davis E. King. “Dlib-ml: A Machine Learning Toolkit”. In: *J. Mach. Learn. Res.* 10 (Dec. 2009). ISSN: 1532-4435. URL: <http://dl.acm.org/citation.cfm?id=1577069.1755843> (p. 60).
- [51] William Koehrsen. *Random Forest Simple Explanation*. URL: <https://medium.com/@williamkoehrsen/random-forest-simple-explanation-377895a60d2d> (p. 33).
- [52] Ying-Fang Lai, Mu-Yen Chen, and Hsiu-Sen Chiang. “Constructing the lie detection system with fuzzy reasoning approach”. In: *Granular Computing* (Nov. 2017). URL: <https://doi.org/10.1007/s41066-017-0064-3> (p. 9).
- [53] Sharon Leal and Aldert Vrij. “Blinking During and After Lying”. In: *Journal of Nonverbal Behavior* 32.4 (Dec. 2008), pp. 187–194. URL: <https://doi.org/10.1007/s10919-008-0051-0> (p. 7).
- [54] Xiaobai Li et al. “A Spontaneous Micro Facial Expression Database: Induce-ment, Collection and Baseline”. In: *Face and Gesture (FG)*. 2013 (p. 14).

- [55] Xiaobai Li et al. “Reading Hidden Emotions: Spontaneous Micro-expression Spotting and Recognition”. In: *IEEE Trans. Affective Computing (TAFFC)*. 2015 (p. 14).
- [56] Kai Keat Lim et al. “Lying Through the Eyes: Detecting Lies Through Eye Movements”. In: *Proceedings of the 6th Workshop on Eye Gaze in Intelligent Human Machine Interaction: Gaze in Multimodal Interaction*. 2013. URL: <http://doi.acm.org/10.1145/2535948.2535954> (p. 8).
- [57] P. Lucey et al. “The Extended Cohn-Kanade Dataset (CK+): A complete dataset for action unit and emotion-specified expression”. In: *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Workshops*. June 2010. DOI: [10.1109/CVPRW.2010.5543262](https://doi.org/10.1109/CVPRW.2010.5543262) (pp. 16, 17, 56, 59, 60).
- [58] Bella M. DePaulo et al. “Lying in Everyday Life”. In: 70 (June 1996), pp. 979–95 (p. 3).
- [59] *Machine Learning 101 what is machine learning*. URL: <https://hydrasky.com/network-security/machine-learning-101-what-is-machine-learning/> (p. 21).
- [60] S. M. Mavadati et al. “DISFA: A Spontaneous Facial Action Intensity Database”. In: *IEEE Transactions on Affective Computing* 4 (Apr. 2013). DOI: [10.1109/T-AFFC.2013.4](https://doi.org/10.1109/T-AFFC.2013.4). URL: doi.ieee.org/10.1109/T-AFFC.2013.4 (pp. 16, 56, 60).
- [61] Carol McDonald. *Demystifying AI, Machine Learning and Deep Learning*. URL: <https://mapr.com/blog/demystifying-ai-ml-dl/> (p. 20).
- [62] G. McKeown et al. “The SEMAINE corpus of emotionally coloured character interactions”. In: *2010 IEEE International Conference on Multimedia and Expo*. July 2010. DOI: [10.1109/ICME.2010.5583006](https://doi.org/10.1109/ICME.2010.5583006) (pp. 16, 56).
- [63] Rada Mihalcea, Verónica Pérez-Rosas, and Mihai Burzo. “Automatic Detection of Deceit in Verbal Communication”. In: *Proceedings of the 15th ACM on International Conference on Multimodal Interaction*. 2013. URL: <http://doi.acm.org/10.1145/2522848.2522888> (p. 6).
- [64] H. Nasri, W. Ouarda, and A. M. Alimi. “ReLiDSS: Novel lie detection system from speech signal”. In: *2016 IEEE/ACS 13th International Conference of Computer Systems and Applications (AICCSA)*. Nov. 2016 (p. 6).
- [65] Machine Learning Group at National Taiwan University. *LIBLINEAR – A Library for Large Linear Classification*. URL: <https://www.csie.ntu.edu.tw/~cjlin/liblinear/> (p. 61).
- [66] Donald Neth and Aleix M. Martinez. “Emotion perception in emotionless face images suggests a norm-based representation”. In: *Journal of Vision* 9.1 (2009). DOI: [10.1167/9.1.5](https://doi.org/10.1167/9.1.5). URL: <http://dx.doi.org/10.1167/9.1.5> (p. 61).
- [67] Annalyn Ng. *Decision Trees: A Disastrous Tutorial*. URL: <https://www.kdnuggets.com/2016/09/decision-trees-disastrous-overview.html> (p. 34).

- [68] D. I. Noje and R. Malutan. “Head movement analysis in lie detection”. In: *2015 Conference Grid, Cloud High Performance Computing in Science (ROLCG)*. Oct. 2015 (p. 12).
- [69] Paola Noreña Cardona. “A COMPENDIUM OF PATTERN RECOGNITION TECHNIQUES IN FACE, SPEECH AND LIE DETECTION”. In: 24 (Nov. 2015) (p. 5).
- [70] M. Owayjan et al. “The design and development of a Lie Detection System using facial micro-expressions”. In: *2012 2nd International Conference on Advances in Computational Tools for Engineering Applications (ACTEA)*. Dec. 2012 (p. 14).
- [71] Verónica Pérez-Rosas et al. “Deception Detection Using Real-life Trial Data”. In: *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction. ICMI ’15*. 2015. URL: <http://doi.acm.org/10.1145/2818346.2820758> (pp. 1, 12, 13, 66, 67).
- [72] Tomas Pfister et al. “Recognising Spontaneous Facial Micro-expressions”. In: *International Conference on Computer Vision (ICCV)*. 2011 (p. 14).
- [73] Stephen M Porter, Leanne ten Brinke, and Brendan B Wallace. “Secrets and Lies: Involuntary Leakage in Deceptive Facial Expressions as a Function of Emotional Intensity”. In: 2012 (p. 3).
- [74] Stephen Porter and Leanne ten Brinke. “The Truth About Lies: What Works in Detecting High-Stakes Deception?” In: 15 (Feb. 2010), pp. 57–75 (p. 5).
- [75] Stephen Porter and Mary Campbell. “A. Vrij, Detecting Lies and Deceit: The Psychology of Lying and Implications for Professional Practice”. In: 7 (Sept. 1999), pp. 227–232 (p. 3).
- [76] Probaway. *What action does a baby’s emotion generate in its mother?* URL: <https://probaway.wordpress.com/tag/dna-of-human-emotional-expression/> (p. 13).
- [77] J. G. Proudfoot et al. “Deception is in the eye of the communicator: Investigating pupil diameter variations in automated deception detection interviews”. In: *2015 IEEE International Conference on Intelligence and Security Informatics (ISI)*. May 2015 (p. 8).
- [78] Quazi Marufur Rahman. *What is a Support Vector Machine?* URL: <https://www.quora.com/What-is-a-Support-Vector-Machine> (p. 43).
- [79] B. Rajoub and R. Zwiggelaar. “Thermal facial analysis for deception detection”. In: *IEEE Transactions on Information Forensics and Security*. 2014 (p. 11).
- [80] Jan Rupnik. *Stochastic Subgradient Approach for Solving Linear Support Vector Machines*. URL: <https://slideplayer.com/slide/8088877/> (p. 44).
- [81] E. Bienenstock S. Geman and R. Doursat. “Neural networks and the bias/variance dilemma”. In: *Neural Computation* 4 (1992) (p. 21).

- [82] Jason M. Saragih, Simon Lucey, and Jeffrey F. Cohn. “Deformable Model Fitting by Regularized Landmark Mean-Shift”. In: *International Journal of Computer Vision* 91.2 (Jan. 2011). ISSN: 1573-1405. DOI: [10.1007/s11263-010-0380-4](https://doi.org/10.1007/s11263-010-0380-4). URL: <https://doi.org/10.1007/s11263-010-0380-4> (p. 54).
- [83] Bj&ouml;rn Schuller et al. “AVEC 2011-the First International Audio/Visual Emotion Challenge”. In: *Proceedings of the 4th International Conference on Affective Computing and Intelligent Interaction - Volume Part II*. Berlin, Heidelberg: Springer-Verlag, 2011. ISBN: 978-3-642-24570-1. URL: <http://dl.acm.org/citation.cfm?id=2062850.2062907> (p. 60).
- [84] Ankit Sharma. *Support Vector Machines without tears*. URL: <https://www.slideshare.net/ankitksharma/svm-37753690> (p. 40).
- [85] A. I. Simbolon et al. “An experiment of lie detection based EEG-P300 classified by SVM algorithm”. In: *2015 International Conference on Automation, Cognitive Science, Optics, Micro Electro-Mechanical System, and Information Technology (ICACOMIT)*. Oct. 2015 (p. 8).
- [86] B. Singh, P. Rajiv, and M. Chandra. “Lie detection using image processing”. In: *2015 International Conference on Advanced Computing and Communication Systems*. Jan. 2015 (p. 8).
- [87] Bruno Stecanella. *An introduction to Support Vector Machines (SVM)*. URL: <https://monkeylearn.com/blog/introduction-to-support-vector-machines-svm/> (p. 41).
- [88] R.G.D Steel and J. H Torrie. *Principles and Procedures of Statistics with Special Reference to the Biological Science*, McGraw Hill. 1960 (p. 29).
- [89] Peter Norvig Stuart J. Russell. *Artificial Intelligence: A Modern Approach, Third Edition - Prentice Hall*. 2010 (p. 20).
- [90] Lin Su and Martin Levine. “Does “lie to me” lie to you? An evaluation of facial clues to high-stakes deception”. In: *Computer Vision and Image Understanding* 147 (2016). URL: <http://www.sciencedirect.com/science/article/pii/S1077314216000345> (pp. 4, 14).
- [91] S. Sumridetchkajorn et al. “Simultaneous Analysis of Far Infrared Signals From Periorbital and Nostril Areas for Nonintrusive Lie Detection: Performance From Real Case Study”. In: *Journal of Lightwave Technology* 33.16 (Aug. 2015) (p. 11).
- [92] *Support Vector Machines*. URL: <http://beta.cambridgespark.com/courses/jpm/05-module.html> (p. 48).
- [93] T. usaitis et al. “OpenFace 2.0: Facial Behavior Analysis Toolkit”. In: *2018 13th IEEE International Conference on Automatic Face Gesture Recognition (FG 2018)*. May 2018 (pp. 50, 55, 58).
- [94] M. F. Valstar et al. “FERA 2015 - second Facial Expression Recognition and Analysis challenge”. In: *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*. May 2015 (p. 60).

- [95] M. F. Valstar et al. “The first facial expression recognition and analysis challenge”. In: *Face and Gesture 2011*. Mar. 2011. DOI: [10.1109/FG.2011.5771374](https://doi.org/10.1109/FG.2011.5771374) (p. 60).
- [96] Bruno Verschuere et al. “The ease of lying”. In: 20 (Nov. 2010), pp. 908–11 (p. 5).
- [97] Aldert Vrij et al. “Police officers ability to detect deception in high stakes situations and in repeated lie detection test”. In: 20 (Sept. 2006), pp. 741–755 (p. 3).
- [98] Wikipedia contributors. *Decision tree learning — Wikipedia, The Free Encyclopedia*. [Online; accessed 4-September-2018]. 2018. URL: https://en.wikipedia.org/w/index.php?title=Decision_tree_learning&oldid=856803699 (p. 34).
- [99] Wikipedia contributors. *Dimensionality reduction — Wikipedia, The Free Encyclopedia*. [Online; accessed 14-September-2018]. 2018. URL: https://en.wikipedia.org/w/index.php?title=Dimensionality_reduction&oldid=846565901 (p. 22).
- [100] Wikipedia contributors. *Functional magnetic resonance imaging — Wikipedia, The Free Encyclopedia*. [Online; accessed 20-June-2018]. 2018. URL: https://en.wikipedia.org/w/index.php?title=Functional_magnetic_resonance_imaging&oldid=842791829 (pp. 9, 10).
- [101] Wikipedia contributors. *Linear regression — Wikipedia, The Free Encyclopedia*. [Online; accessed 5-August-2018]. 2018. URL: https://en.wikipedia.org/w/index.php?title=Linear_regression&oldid=851432903 (pp. 27, 28).
- [102] Wikipedia contributors. *Logistic regression — Wikipedia, The Free Encyclopedia*. [Online; accessed 19-August-2018]. 2018. URL: https://en.wikipedia.org/w/index.php?title=Logistic_regression&oldid=854182949 (p. 30).
- [103] Wikipedia contributors. *Mel-frequency cepstrum — Wikipedia, The Free Encyclopedia*. 2018. URL: https://en.wikipedia.org/w/index.php?title=Mel-frequency_cepstrum&oldid=835415759 (p. 5).
- [104] Wikipedia contributors. *Point distribution model — Wikipedia, The Free Encyclopedia*. [Online; accessed 9-August-2018]. 2018. URL: https://en.wikipedia.org/w/index.php?title=Point_distribution_model&oldid=832628847 (pp. 53, 54).
- [105] Wikipedia contributors. *Random forest — Wikipedia, The Free Encyclopedia*. 2018. URL: https://en.wikipedia.org/w/index.php?title=Random_forest&oldid=854167685 (pp. 33, 36).
- [106] Wikipedia contributors. *Regression analysis — Wikipedia, The Free Encyclopedia*. [Online; accessed 3-August-2018]. 2018. URL: https://en.wikipedia.org/w/index.php?title=Regression_analysis&oldid=850022021 (pp. 26, 27).

- [107] Wikipedia contributors. *Statistical classification — Wikipedia, The Free Encyclopedia*. [Online; accessed 4-September-2018]. 2018. URL: https://en.wikipedia.org/w/index.php?title=Statistical_classification&oldid=856605904 (p. 24).
- [108] Zhe Wu et al. “Deception Detection in Videos”. In: *CoRR* abs/1712.04415 (2017). URL: <http://arxiv.org/abs/1712.04415> (p. 12).
- [109] A. Zadeh, T. Baltrušaitis, and L. Morency. “Convolutional Experts Constrained Local Model for Facial Landmark Detection”. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. July 2017. DOI: [10.1109/CVPRW.2017.256](https://doi.org/10.1109/CVPRW.2017.256) (pp. 51, 53).
- [110] Xing Zhang et al. “BP4D-Spontaneous: a high-resolution spontaneous 3D dynamic facial expression database”. In: *Image and Vision Computing* 32.10 (2014). DOI: <https://doi.org/10.1016/j.imavis.2014.06.002>. URL: <http://www.sciencedirect.com/science/article/pii/S0262885614001012> (pp. 16, 17, 56).

List of Figures

1.1	Performance on deceit detection by human observers [90].	4
1.2	Topside is the spoken word, bottom is the MFCC derived from the word	6
1.3	Eyes could hold a lot of informations regarding what we are thinking [27].	7
1.4	EEG of P300 waves image on channels Fz, Cz, and Pz	9
1.5	fMRI image with yellow areas showing increased activity compared with a control condition [100]	10
1.6	Examples of thermal images during (a) questioning and (b) answering [91]	11
1.7	Distribution of non-verbal features for deceptive and truthful groups [71].	13
1.8	Six basic Facial Expressions in adults and children [76].	13
1.9	Duchenne smile (real) vs non Duchenne (fake).	15
1.10	Facial images with AU intensity variations (AU12: lip corner puller, AU25: lips part) [60].	16
1.11	AUs of six compound facial expressions of emotion. The AUs of the basic emotions are combined to produce the compound category [29].	16
2.1	Machine Learning Subfields [61]	20
2.2	Supervised Learning Workflow [59]	21
2.3	Example of Bias-Variance Tradeoff [17]	22
2.4	Example of overfitting and underfitting [15]	23
2.5	Example of Unclustered and Clustered data [48]	24
2.6	Example of data division by a classification algorithm	25
2.7	Regression Analysis can solve this! [10]	26
2.8	Example of Regression Analysis	27
2.9	The fit is found by minimizing the sum of squared errors. Each gray line segment represents an error, and the prediction makes a compromise by averaging their square [37]	28
2.10	Graph of a logistic regression curve showing probability of passing an exam versus hours studying [102]	30
2.11	Example of a Sigmoid shaped function	31
2.12	Example of Random Forest Classification [51]	33
2.13	Example of a Decision Tree [67]	34
2.14	Bagging Example [1]	36
2.15	Bagging leads to different decision trees [2]	37

2.16 Random Forest overview	38
2.17 Example of variable importance found with Random Forest [40].	39
2.18 Two and three dimensional space hyperplanes [84]	40
2.19 Different hyperplanes have different margins [87]	41
2.20 Example of Support Vectors	42
2.21 Example of Linearly non separable data [78]	43
2.22 Adding a single data point (from left to right) can change the maximal margin hyperplane very significantly [37]	44
2.23 Soft margin classifier. Some data points are allowed to be misclassified to find a larger margin [80]	44
2.24 Visualization of different kinds of kernels [34].	47
2.25 Linearly non separable data [92].	48
2.26 Transformation of the data in a feature space where the instances from the two classes may be linearly separable [92].	48
3.1 OpenFace Pipeline [93]	50
3.2 Facial landmarks naturally cluster around appearance prototypes (facial hair, expressions, make-up etc). To model such appearance variations a Convolutional Experts Network is used to merge the advantages of neural architectures and mixtures of patch experts, to model landmark alignment probability [109].	51
3.3 Overview of CLM [25].	52
3.4 Overview of the Convolutional Experts Network model. The output response map is a non-negative and non-linear combination of neurons in ME-layer using a sigmoid activation [109].	53
3.5 PDM of a metacarpal. Dots mark the possible position of landmarks, and the line denote the mean shape [47].	53
3.6 Example of different facial landmarks detected in diverse conditions and viewing angles. [93].	55
3.7 AU detection and intensity pipeline [12].	56
3.8 AUs coded in the DISFA database [28].	57
3.9 List of available AUs for prediction [93].	58
3.10 Stable points used for alignment to a common reference frame, followed by masking [12].	60
3.11 Example of input data in csv format.	62
4.1 Examples of images from the dataset videos. [71].	66
4.2 Frequency of gestures annotated in the videos from the database. [71].	67
4.3 Average AU Occurrences for the deceptive training set.	69
4.4 Average AU Occurrences for the truthful training set.	70
4.5 Comparison of AU Occurrences for the training set.	70
4.6 Correlation between variables for both presence and intensity in the training set.	71
4.7 p-values from GLM for each AU (intensity and presence).	72
4.8 Variable Importance extracted from classifying data with Random Forest.	73