

Contents

2	Background and Literature Review	2
2.1	Problem Space for Facial Expression Analysis	2
2.1.1	Level of Description	2
2.2	System Structure	4
2.3	Face Detection	5
2.3.1	Eigenface and Template Matching	6
2.4	Curve Fitting	8
2.5	Feature Extraction	9
2.6	Tensorface	10
2.7	Potential Net	11
2.7.1	Support Vector Machine	12
2.8	Face Database	14
2.8.1	Databases For Identity Recognition	14
2.8.2	Databases for Expression Recognition	17
2.9	Chapter Summary	18

Chapter 2

Background and Literature Review

The importance of facial expression in social interaction and social intelligence is widely recognized. Facial expression analysis has been an active research topic since 19th century. The first automatic facial expression recognition system was introduced in 1978 by Suwa et al. [83]. This system attempts to analyze facial expressions by tracking the motion of 20 identified spots on an image sequence. Since then, a lot of work has been done in this domain. Various computer systems have been made to help us understand and use this natural form of human communication.

This chapter reviews the state of the art of what has been done in processing and understanding facial expression. When building an FER system, these main issues must be considered: face detection and alignment, image normalization, feature extraction, and classification. Most of the current work in FER is based on methods that implement these steps sequentially and independently. Before exploring what has been done in literature for implementing these steps, we will briefly describe the problem space for facial expression analysis.

2.1 Problem Space for Facial Expression Analysis

2.1.1 Level of Description

In general there are two types of method to describe facial expression.

Facial Action Coding System

The facial action coding system [24] is a human-observer-based system widely used in psychology to describe subtle changes in facial features. FACS consists of 44 action units which are related to contraction of a specific set of facial muscles (Fig.2.1). Some of the action units are shown in Fig.2.2. Conventional, FACS code is manually labeled by trained observers while viewing videotaped facial behavior in slow motion. In recent years, some attempts have been made to do this automatically [69]. The advantage of FACS is its ability to capture the subtlety of facial expression, however FACS itself is purely descriptive and includes no inferential labels. That means in order to get the emotion estimation, the FACS code needs to be converted into the Emotional Facial Action System (EMFACS [28]) or similar systems.

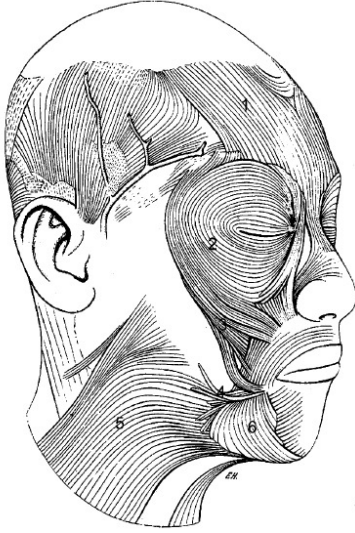


Figure 2.1: Muscles of facial expression. 1, frontalis; 2, orbicularis oculi; 3, zygomaticus major; 4, risorius; 5, platysma; 6, depressor anguli oris [33]

Upper Face Action Units					
AU 1	AU 2	AU 4	AU 5	AU 6	AU 7
Inner Brow Raiser	Outer Brow Raiser	Brow Lowerer	Upper Lid Raiser	Cheek Raiser	Lid Tightener
*AU 41	*AU 42	*AU 43	AU 44	AU 45	AU 46
Lid Droop	Slit	Eyes Closed	Squint	Blink	Wink
Lower Face Action Units					
AU 9	AU 10	AU 11	AU 12	AU 13	AU 14
Nose Wrinkler	Upper Lip Raiser	Nasolabial Deepener	Lip Corner Puller	Cheek Puffer	Dimpler
AU 15	AU 16	AU 17	AU 18	AU 20	AU 22
Lip Corner Depressor	Lower Lip Depressor	Chin Raiser	Lip Puckerer	Lip Stretcher	Lip Funneler
AU 23	AU 24	*AU 25	*AU 26	*AU 27	AU 28
Lip Tightener	Lip Pressor	Lips Part	Jaw Drop	Mouth Stretch	Lip Suck

Figure 2.2: FACS action units [35]

Prototypic Emotional Expressions

Instead of describing the detailed facial features, most FER systems attempt to recognize a small set of prototypic emotional expressions. The most widely-used set is perhaps human universal facial expressions of emotion which consists of six basic expression categories that have been shown to be recognizable across cultures 2.3 .

These expressions, or facial configurations have been recognized in people from widely divergent cultural and social backgrounds, and they have been observed even in the faces of individuals born deaf and blind.

These 6 basic emotions, *i.e.*, disgust, fear, joy, surprise, sadness and anger plus “neutral” which means no facial expression are considered in this work. Given a facial image, our system either works as a conventional classifier to determine the most likely emotion or estimates the weights (or possibility) of each emotion as a fuzzy classifier does.

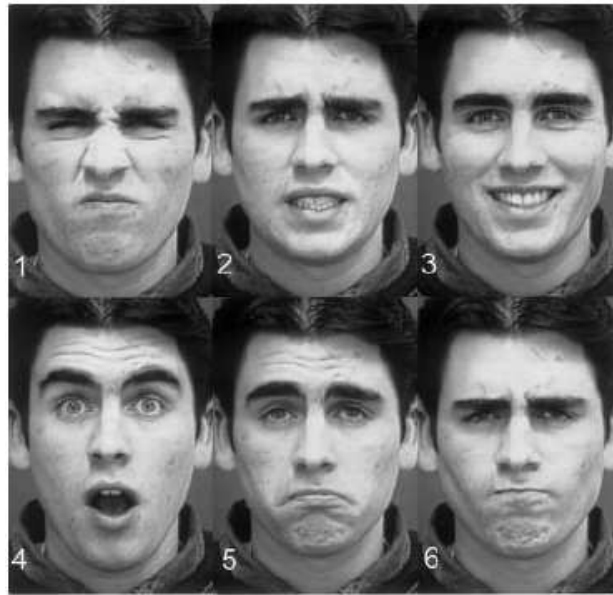


Figure 2.3: Basic facial expression phenotypes. 1, disgust; 2, fear; 3, joy; 4, surprise; 5, sadness; 6, anger

2.2 System Structure

FER can be considered as a special face recognition system or a module of a face recognition system. So it should be instructive to look at the general architecture of a face recognition system. Normally, it consists of four components as depicted in 2.4

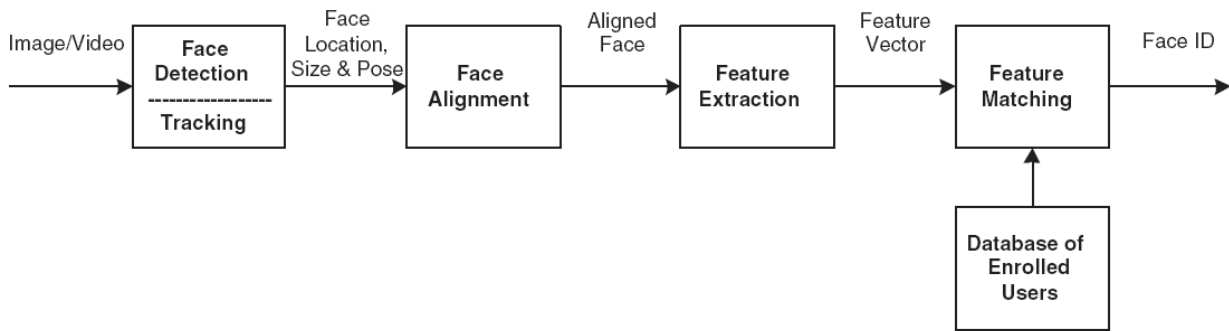


Figure 2.4: Face recognition processing flow

Face detection finds the face areas in the input image. If the input is a video, to be more efficient and also to achieve better robustness, face detection is only performed on key frames and a tracking algorithm is applied on interval frames. Face alignment is very similar to detection, but it is aimed at achieving a more accurate localization. In this step, a set of facial landmarks (facial components), such as eyes, brows and nose, or the facial contour are located; based on that, the face image is rotated, chopped, resized and even warped, this is called geometrical normalization. Usually the face is further normalized with respect to photometrical properties such as illumination and gray scale.

Feature extraction is performed on a normalized face to provide effective information that should be useful for recognizing and classifying labels in which there is interest, such as identity, gender, or expression. The extracted feature vector is sent to a classifier and compared with the training data to produce a recognition output.

2.3 Face Detection

Face detection is the first step in face recognition. It has a major influence on the performance of the entire system. Several cues can be used for face detection, for example, skin color, motion (for videos), facial/head shape, and facial appearance. Most successful face detection algorithms are based on only appearance. This may be because

appearance-based algorithms avoid difficulties in modeling 3D structures of faces. However, the variations of 3D structures due to facial expression and head pose actually heavily affect the facial appearance and make the face/non-face boundary highly complex [7]. To deal with this, a vast arrange of methods have been proposed since the 1990s.

Turk and Pentland [87] describe a detection system based on eigen decomposition which is also known as principal component analysis (PCA). In their method, an image is represented by an average face plus a set of weighted eigenfaces. Whereas only the face images are considered in eigenface, Sung and Poggio [82] also consider the distribution of non-face images and apply Bayes rule to obtain a likelihood estimation. Rowley et al. [72] use neural networks and Osuna et al. [68] trained a Kernel Support Vector Machine to classify face and non-face images. In these systems, a bootstrap algorithm is iteratively used to collect meaningful examples for retraining the detector.

Schneiderman and Kanade [75] use AdaBoost learning to construct a classifier based on wavelet representation of the image. This method is computationally expensive because of the wavelet transformation. To overcome this problem, Viola and Jones [93] replace wavelets with Haar features, which can be computed very efficiently [20] [80]. Their system is the first realtime frontal-view face detector [94].

Under Violas framework, some improvements have been proposed. Lienhart et al. [52] use rotated Haar features to deal with in-plane rotation. Li et al. [51] [50] [52] propose a multiview face detection system which can also handle out-of-plane rotation using a detector-pyramid.

In the following sections, we will describe two face detection algorithms: Eigenface is one of the simplest methods and Violas framework may be the most successful one. AdaBoost learning is an important component in Violas framework and this algorithm will also be useful in the feature extraction module, so our presentation focuses on this part.

2.3.1 Eigenface and Template Matching

Eigenface assumes the face image $x = (x_1, x_2, \dots, x_N)$ is amenable to a multivariate normal distribution from which the training images are identically independently drawn.

This distribution can be described by the following probability density function:

$$f(x_1, x_2, \dots, x_N) = \frac{1}{(2\pi)^{N/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu))$$

where Σ is the covariance matrix and μ is the expectation of x . Eigenface decomposes Σ using eigen decomposition as

$$\Sigma = USU^T$$

where U is a unitary matrix and $S = \text{diag}(s_1^2, s_2^2, \dots, s_N^2)$ is a diagonal matrix with all elements non-negative. Each column of U , U_i , is called an Eigenface. A face image x can be represented by μ and U_i as

$$x = \mu + \sum_i a_i U_i$$

It can be shown that $\frac{a_i}{s_i}$ are i.i.d. standard normal variables. So the probability density function of x is:

$$f(x) = \prod_i \frac{1}{(2\pi)^{1/2}} \exp(-\frac{1}{2} \frac{a_i^2}{s_i^2}) = \frac{1}{(2\pi)^{N/2}} \exp(-\frac{1}{2} \sum_i \frac{a_i^2}{s_i^2})$$

Equation (2.4) can be used as a probability estimate and we can define a distance measure according to

$$D = \sum_i \frac{a_i^2}{s_i^2}$$

which is called normalized Euclidean distance. A large D implies a small probability of being a face image and vice versa. Based on (2.5) Turk and Pentland [87] built a face detection system. Sung and Poggio's paper [82] used a similar idea, they assume images are produced by a mixture of Gaussian models: a face image Gaussian and a non-face image Gaussian. So they also estimate the probability of x of being a non-face image $f'(x)$ and the final decision is made using a Bayesian classifier.

If we further assume Σ is an identity matrix, (2.5) degenerates into Euclidean distance which means the probability density function is controlled by, $|x - \mu|^2$, the variation of the image from the average face μ . This gives the simplest detection algorithm: template matching, i.e., finding a "face template" μ , and then for each x determine whether it is a face image by thresholding $|x - \mu|$.

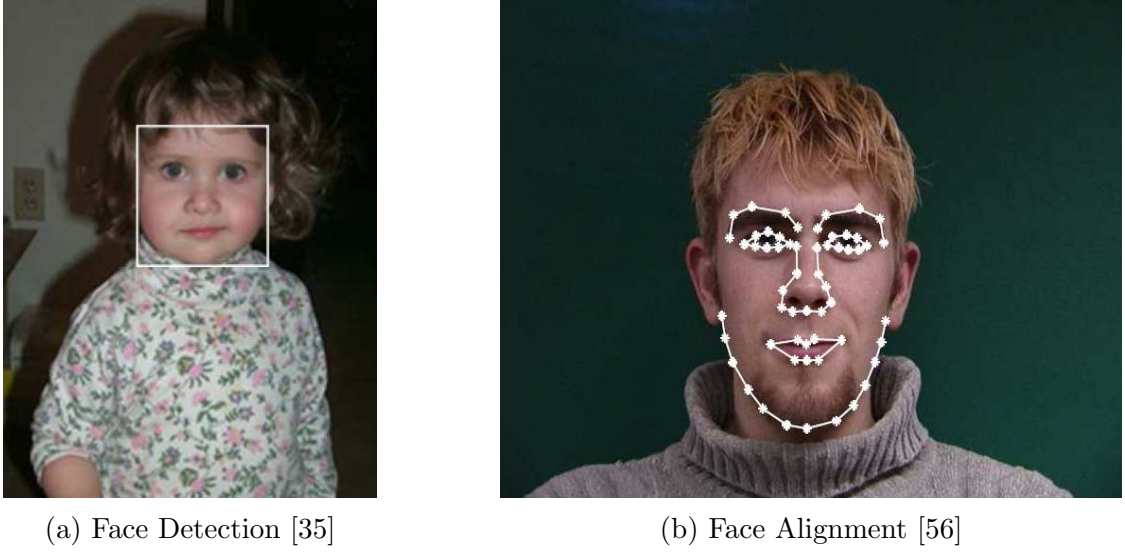


Figure 2.5: Detection Vs Alignment

2.4 Curve Fitting

The basic problem of curve fitting is to locate the contour of an object in an image. Most curve fitting methods (explicitly or implicitly) consider two forces: internal force (elastic force) caused by deformation and external force (image force) caused by density gradient. As in physics, both internal force and external force cause potential energy, and the curve fitting is achieved by minimizing the overall energy functional (Fig.2.8).

In general, a contour can be represented by $c = c(s)$, parameterized by its arc length, s . The energy function can be expressed by:

$$\varepsilon = E_{elastic}(c) + E_{image}(c) = \int_c (a(s)e_{elastic}(c, s) + b(s)e_{image}(c, s))ds \quad (2.1)$$

So the optimization problem can be expressed as

$$c \arg \min E_{elastic} + E_{image}(c) \quad (2.2)$$

Usually, (2.24) is solved using iterative searching algorithms. A large number of curve fitting methods have been proposed using a wide range of energy functions and searching schemes.

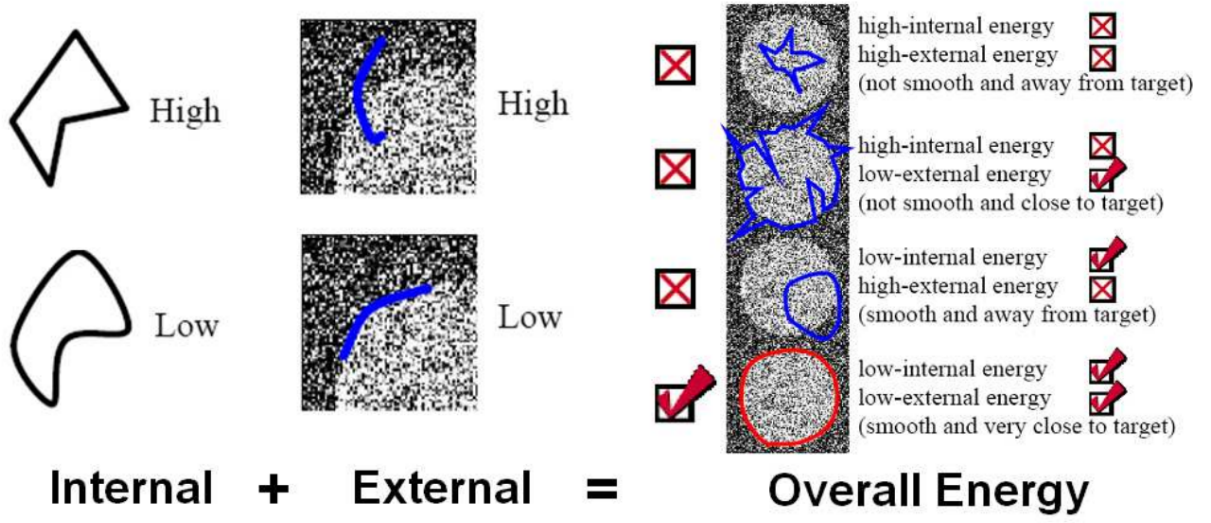


Figure 2.6: Energy function and curve fitting

It's worth pointing out that sometimes the energy function is not explicitly given, instead it's implicitly defined in the searching scheme, for example, in Active Shape Model.

2.5 Feature Extraction

Feature extraction converts pixel data into a higher-level representation of shape, motion, color, texture, and spatial configuration of the face or its components. The extracted representation is used for subsequent classification. Feature extraction generally reduces the dimensionality of the input space. The reduction procedure should (ideally) retain essential information possessing high discrimination power and high stability [14]. In the face recognition area, various features have been used.

The coefficients of Eigenface can be used as features and recently, an extension of Eigenface [90] defines Tensorface which has shown a promising choice of feature. Active Appearance Model [17] decomposes the facial image into shape and texture. The shape vector which is coded using ASM describes the contour of the facial components, whereas the texture vector gives the shape-free facial texture. Matsuno et al. [41] extract features

using a two-dimensional mesh, called Potential Net. All the above-mentioned methods are considered as holistic features, because they are related to the overall structure of the image. There is another kind of features called local features, each of which focuses only on a small region. The most straightforward idea may be to directly use image sub-windows as local features: for example, in [?] Colmenarez *et al.* use nine sub-windows located around the facial components. Wavelet filters have been used too, the most popular of which is the Gabor filter which has been shown [?] [?] to be a reasonable model of visual processing in primary visual cortex. Yin and Wei use topographic primitive features to represent faces [?]. In [?], instead of defining features ahead of time, Yu and Bhanu use a evolutionary algorithm to generate features automatically. For video-based FER, the dynamic of expression can also serve as features. [?] proposes Geometric Deformation Feature which represents the geometrical displacement of certain selected landmark nodes. In [?], Aleksic and Katsaggelos use Facial Animation Parameters which are based on Active Shape Model.

In this section, we'll briefly introduce some of these feature extraction methods

2.6 Tensorface

Tensorface [?] is a multilinear extension of Eigenface. Instead of representing the face image using a linear equation.

$$x = \bar{x} + \sum_i a_i x_i \quad (2.3)$$

It models the face by a multilinear system which is equivalent to

$$x = \bar{x} + \sum_{i_1} \dots \sum_{i_N} a_{i_N} \dots a_{i_1} x_{i_1 \dots i_N} \quad (2.4)$$

Compared to Eigenface which ignores the label of images, Tensorface analyzes a face ensemble with respect to its underlying factors(labels): for example, identities, views, and illuminations. The "principal components" in this multilinear system are referred to as Tensorfaces which are shown in Fig.2.9. The Tensorface coefficients a_{i_k} can be used as features in a recognition task, and because the original image can be reproduced using 2.26, Tensorface coefficients can also be used for image synthesis where we first generate a set of coefficients, then use them to synthesize images.

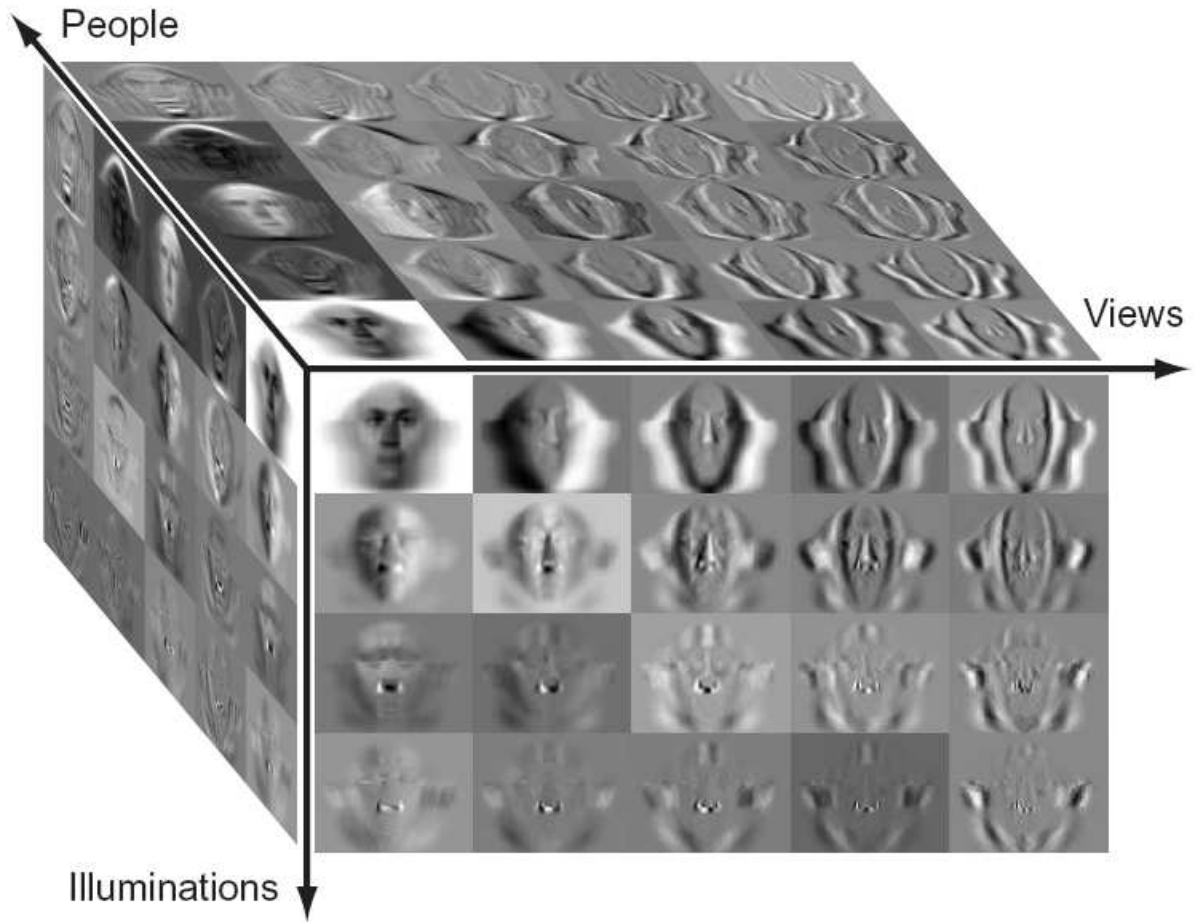


Figure 2.7: A partial visualization of TensorFaces bases for an ensemble of 2,700 facial images spanning 75 people, each imaged under 6 viewing and 6 illumination conditions [?]

Since Tensorface is shown to be a promising method in face recognition, some improvements have been proposed. [?] proposes Multilinear Independent Component Analysis where they try to find the independent directions of variation. In [?] Shashua et al. introduce Non-Negative Tensor Factorization which is a generalization of Non-negative Matrix Factorization.

2.7 Potential Net

Matsuno *et al.* [?] [?] propose Potential Net to extract facial features. As shown in Fig.2.10, Potential Net is a two dimensional mesh of which nodes are connected to their four neighbors with springs, while the most exterior nodes are fixed to the frame of the Net. Similar to curve fitting, Potential Net considers two forces: each node in the mesh is driven

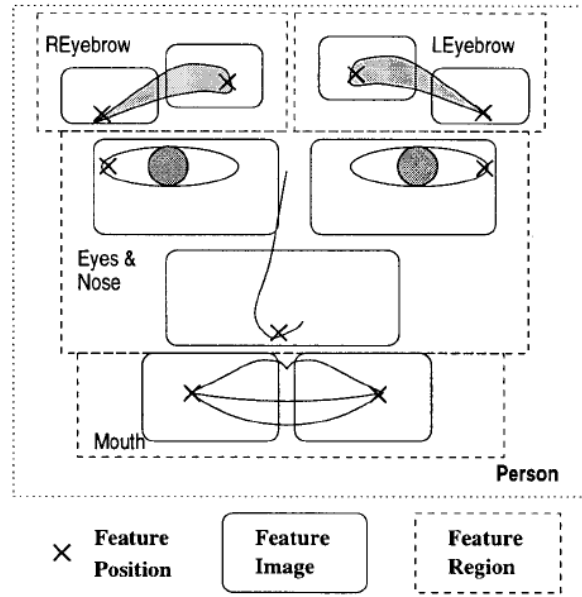


Figure 2.8: Figure 2.12: Scheme of the Facial Features and Regions [15]

2.7.1 Support Vector Machine

Support Vector Machine attempts to construct a linear classifier which maximizes the margin between two classes, so its also known as Optimal Margin Classifier [9]. Fig.2.13 gives an SVM classifier where $\frac{1}{|w|}$ gives the margin and samples along the hyper-planes are called the support vectors.

It has been proven that SVM minimizes the Structural Risk Function which is considered as a better error estimation than the normallyused Empirical Risk Function in terms of generalization capacity.

We consider data points of the form: $\{(x_1, y_1), \dots, (x_n, y_n)\}$ where y_i is either 1 or -1, a label denoting the class to which the point x_i belongs. The basic version of SVM can be written as

$$\begin{aligned} \text{argmax} \quad & \frac{1}{||w||} \\ \text{s.t.} \quad & y_i(x^T w - b_0) \geq 1, \quad \forall i \end{aligned}$$

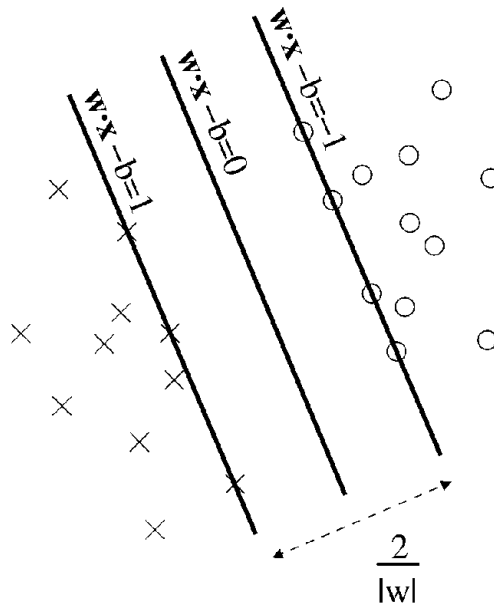


Figure 2.9: Figure 2.13: Maximum-margin hyper-planes for a SVM trained with samples from two classes [99]

$||\dots||$ in (2.27) can be replaced by any distance measure. If norm-2 is used, the problem is equivalent to

$$\begin{aligned} \operatorname{argmin} \quad & \frac{1}{2} ||w||^2 \\ \text{s.t.} \quad & y_i(x^T w - b_0) \geq 1, \quad \forall i \end{aligned}$$

Equation (2.28) is a quadratic programming and according to the strong duality theorem it can be converted to:

$$\begin{aligned} \operatorname{argmax} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad \forall i \end{aligned}$$

which is called the dual problem of (2.28). In practice, we always use (2.29) as it is easier to handle numerically. Moreover, in (2.29) all the computation of x_i is written in terms of inner product, and that means it can be generalized to a nonlinear case by employing Kernel technique.

2.8 Face Database

”Because of its non rigidity and complex three-dimensional structure, the appearance of a face is affected by a large number of factors including identity, face pose, illumination, facial expression, age, occlusion, and facial hair. The development of algorithms robust to these variations requires databases of sufficient size that include carefully controlled variations of these factors. Furthermore, common databases are necessary to comparatively evaluate algorithms. Collecting a high quality database is a resource-intensive task: but the availability of public face databases is important for the advancement of the field” [35]. In this section we briefly review some publicly available databases for face recognition, face detection, and facial expression analysis, and we’ll mainly focus on the three databases which we will use in this thesis.

To facilitate this statement, we divide face databases into two categories according to their designing goals. In the first part, we’ll introduce databases which are normally used for face recognition; those which are dedicated to expression recognition will be discussed in the second part. As only a few databases are of the second type, and FER system shares some common modules with identity recognition system, in this work we also use some databases of the first type.

2.8.1 Databases For Identity Recognition

Most face databases are of this category (Table.2.1). To test for robustness, some of them are captured under different poses, illuminations and expressions. However, because they’re mainly designed for identity recognition, the expressions are added as noise and usually not well controlled. So in general these databases are considered not suitable for FER research. In our work, we only use them to train peripheral modules (processing and Feature Extraction).

The IMM Face Database [66]

The IMM Face Database comprises 240 still images of 40 individuals (7 females and 33 males), all without glasses. For each person, 6 images are provided:

Table 2.1: Some of the most popular Face Recognition Databases [35]

Database	No. of subjects	Pose	Illumination	Facial Expressions
AR	116	1	4	4
BANCA	208	1	++	1
CAS-PEAL	66-1040	21	9-15	6
CMU HYPER	54	1	4	1
CMU PIE	54	1	4	1
Equinox IR	91	1	3	3
FERET	1199	9-20	2	2
Harvard RL	10	1	77-84	1
IMM FACE	40	3	2	3+
KFDB	1000	7	16	5
MIT	15	3	3	1
MPI	200	3	3	1
ND HID	300+	1	3	2
NIST MID	1573	2	1	++
ORL	10	1	++	++
UMIST	20	++	1	++
U.Texas	284	++	1	++
U. Oulu	125	1	16	1
XM2VTS	295	++	1	++
Yale	15	1	3	6
Yale B	10	9	64	1

- Frontal face, neutral expression, diffuse light.
- Frontal face, happy expression, diffuse light.
- Face rotated approx. 30 degrees to the persons right, neutral expression, diffuse light.
- Face rotated approx. 30 degrees to the persons left, neutral expression, diffuse light.

- Frontal face, neutral expression, spot light added at the persons left side.
- Frontal face, joker image (arbitrary expression), diffuse light.

The images are stored in 640 × 480 JPEG files. Owing to technique problems, most images are RGB, but some are grey-scale [66]. One good thing about this database is that manually labeled face contour is available. The following facial structures were annotated using 58 landmarks: eyebrows, eyes, nose, mouth and jaw. These landmarks are divided into seven point paths; three closed and four open as shown in Fig.2.14. In our work, this database will be used to train the ASM and AAM model.



Figure 2.10: Example image from IMM face database

CMU Pose, Illumination, and Expression Database [79]

The CMU-PIE database is among the most comprehensive databases in this area. It systematically samples a large number of pose and illumination conditions along with a variety of facial expressions. The PIE database was captured under 21 illuminations (lit by 21 flashes) from 13 directions (using 13 synchronized cameras). In total, there are 41,368 images obtained from 68 individuals. In our experiment, we only use a sub-set of this database which consists of images of 62 people. 25 images were selected for each individual

with 5 different viewpoints and 5 different illuminations. Part of the data set is shown in Fig.2.15.

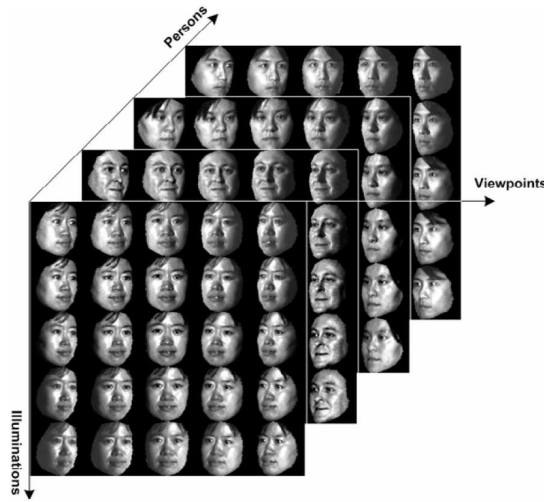


Figure 2.11: A subset of CMU PIE database [53]

2.8.2 Databases for Expression Recognition

”The human face is able to display an astonishing variety of expressions. Collecting a database that samples this space in a meaningful way is a difficult task” [35]. As a result, there are many fewer databases available for expression recognition (Table 6.1). As mentioned in 2.1.1, there are two ways to describe facial expressions. Available databases can be categorized into two classes according to the description they used. In one group [38] expressions are coded in FACS, while in the other group [57] images are labeled by their prototypic emotional expressions.

Japanese Female Facial Expression Database [57]

The JAFFE database contains 213 images of 10 Japanese female models. Their images are labeled by emotions: six basic emotions (anger, disgust, fear, joy, happy, sad and surprise) are considered and Neutral is added as the 7th emotion which is defined through the absence of expression. Fig.2.16 shows example images for one subject along with emotion

Table 2.2: Commonly used expression recognition databases [35]

Database	No. of subjects	No. of Expressions	Image Resolution	Video/Image
JAFFE	10	7	256 X 256	Image
U. Maryland	40	6	560 X 240	Video
Cohn-Kanade	100	23	640 X 480	Video



Figure 2.12: Example images from JAFFE database[35]

labels. The images were originally printed in monochrome and then digitized using a flatbed scanner.

2.9 Chapter Summary

In this chapter, we first talked about the background of facial analysis, then gave an overview of the development in this area, and we also briefly introduced some state of the art techniques which might be useful for our system. At the end, we had a glance at some face databases for identity and expression recognition. Starting in the next chapter, we'll discuss the design of our FER system.