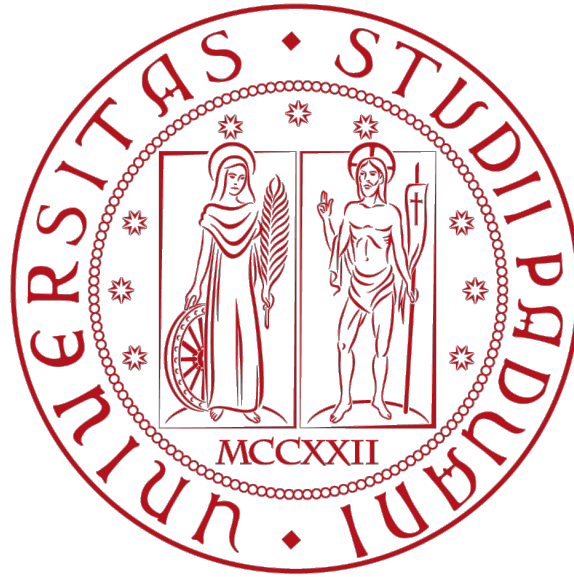


**UNIVERSITÀ DEGLI STUDI DI PADOVA**

Dipartimento di ingegneria dell'informazione  
Corso di Laurea in Ingegneria dell'Automazione

**IMPERIAL COLLEGE OF LONDON**

Department of Computing



**Imperial College  
London**

**Dynamical PLDA for face recognition in videos**

*Laureando:*

Alessandro  
FABRIS

*Relatori:*

Prof. Ruggero  
CARLI

Prof. Stefanos  
ZAFEIRIOU

A.A. 2014/2015

Padova, 06/10/2015



## Abstract

Face recognition is a thriving area of computer vision, with applications spanning the areas of home entertainment, information security, law enforcement and surveillance. In this context, PLDA is a well established generative model that guarantees good discriminability between different individuals and great flexibility, hence being suitable for many different tasks. Within this document we provide a brief primer about face recognition in videos, outlining challenges, applications, different approaches and peculiarities with respect to the domain of still images. We then recall PLDA models, along with their founding ideas and the mathematics necessary to exploit them. Given these premises, we present the first dynamical generalization of PLDA (DPLDA), closely related to Kalman filters, which naturally apply to video sequences. We anticipate a performance improvement for DPLDA in comparison with classical PLDA, providing a theoretical justification in the area of face recognition in videos. Subsequently, we present a robust recognition pipeline, suitable for identity inference in videos under uncontrolled pose, illumination, expression, occlusion and duration. We apply this pipeline to challenging recognition tasks for in-the-wild databases, showing good performance with low- and medium-quality videos of less than 3 seconds. Furthermore, we compare PLDA and DPLDA within this pipeline, showing how DPLDA systematically outperforms PLDA. Finally, we outline future research directions to build on this work and further improve recognition performance.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Applications . . . . .	2
1.2	Peculiarities and challenges . . . . .	3
1.3	Approaches . . . . .	3
<b>2</b>	<b>PLDA</b>	<b>11</b>
2.1	A preliminary review of LDA . . . . .	11
2.2	PLDA: a static model for inference about identity . . . . .	11
2.3	Recognition . . . . .	14
2.4	Learning . . . . .	16
2.4.1	E step . . . . .	17
2.4.2	M step . . . . .	18
<b>3</b>	<b>Dynamical PLDA</b>	<b>21</b>
3.1	PLDA for videos . . . . .	21
3.2	Learning and recognition . . . . .	23
3.2.1	M step . . . . .	24
3.2.2	E step . . . . .	25
3.2.3	Recognition . . . . .	28
3.3	Implementation Details . . . . .	28
3.3.1	Video concatenation . . . . .	28
3.3.2	Offline computations . . . . .	30
3.3.3	Inference . . . . .	30
<b>4</b>	<b>Pipeline and simulations</b>	<b>33</b>
4.1	Face frontalization . . . . .	33
4.2	Feature extraction . . . . .	34
4.3	Simulations . . . . .	36
4.3.1	Synthetic data . . . . .	36
4.3.2	Real data . . . . .	37
<b>5</b>	<b>Conclusion</b>	<b>41</b>
	<b>Appendices</b>	<b>43</b>
<b>A</b>	<b>The EM Algorithm</b>	<b>43</b>
	<b>Acknowledgments</b>	<b>45</b>
	<b>Bibliography</b>	<b>47</b>



# 1. Introduction

Individual recognition is a fundamental ability for the members of an organized society. The possibility for humans to distinguish between different people is a self-evident prerequisite for the establishment of meaningful interpersonal relationships, which we reasonably overlook and take for granted. At a more or less subconscious level, the human brain has employed face recognition and biometrics since mankind itself exists, most likely with satisfactory results for everyday life applications. However, as human communities grew, purely natural skills became less suitable for recognition tasks, as their scope and difficulty grew proportionally to populations.

Human-to-human recognition is particularly faulty on large-scale applications, wherein a systematic approach is usually beneficial. Especially during the 19th century, cities grew larger and their population changed quicker than ever thanks to migration flows from the countryside and more accessible forms of transportation. Due to these changes, the location-specific knowledge, law enforcement officers relied on to keep crime under control, lost part of its effectiveness. At the same time the institution of recidivism was introduced in most jurisdictional European systems as an aggravating factor for law violations. These facts concurred in determining the need to establish a systematic way of recording law infringement in conjunction with the identity traits of the offender.

In the 1870s, pioneering French criminologist officer Alphonse Bertillon established a number of combined biometrics techniques – later named “Bertillonage” – that categorized criminals based on available pictures, descriptions and body measurements from a sliding compass. Furthermore he was the first to develop the two-part “mug-shot”, consisting of a frontal and a profile picture against a measuring background, still widely used today. Two decades later, in the 1890s, Argentine police official Juan Vucetich successfully introduced fingerprints as a means of identification, which quickly became an effective technique employed all over the world. In addition, in 1936, ophthalmologist Franck Burch proposed to use iris patterns as a means to identify individuals.

It wasn’t until the emergence of computers however, that the first automatic face recognition system was deployed by Woodrow Bledsoe and colleagues in 1964. Research and development was mainly funded by governmental agencies since applications such as home entertainment were not in sight yet. A fundamental breakthrough took place in 1987 when Sirovich and Kirby [26] applied Principal Component Analysis to face recognition and showed that less than 100 features properly extracted from a face image can be enough to retain most of its information content, thus opening the way for scalable applications to come in the consequent years.

Nowadays biometrics is ubiquitous and relies on both physiological traits, such as voice timbre and rhythm, and behavioural features such as gait and typing patterns. Face appearance is a fundamental part of biometrics and possibly the one with most applications, spanning the areas of home entertainment, information security, law enforcement and surveillance. This technology is now mature after several decades of research, thus

---

achieving good performance in terms of recognition accuracy and speed. At the same time it has the inherent advantage of not requiring the collaboration of the individual to be identified – as opposed to fingerprints or iris pattern acquisition – guaranteeing a good trade off between performance and obtrusiveness. Nevertheless permanence is a critical factor, in that face appearance is bound to change over time, thus requiring a continuous and expensive update of databases. Research about face ageing has been carried out (see [23] for a good review), however its impact on the performance of face recognition systems is still widely unknown. For these reasons, face identification is regarded as a mature technology, suitable for many applications, yet still widely researched and susceptible to improvements.

Broadly speaking, face recognition can be divided into 3 categories, depending on the face data available: algorithms based on still image data, those that exploit video sequences with temporal information and those that rely on more elaborate and higher-dimensional information, possibly comprehensive of infrared-imagery and depth maps. The rest of this chapter provides a brief introduction to face recognition in video sequences, highlighting peculiarities, challenges, approaches and applications. For a complete and fairly recent survey about the topic see [3].

## 1.1 Applications

Up to this point, we have introduced identity inference tasks in generic terms, using the words recognition and identification more or less interchangeably. In this section we briefly characterize four different applications, highlighting their differences.

Face **verification** is a task of authentication whereby a subject claims to be someone. Hence, given the individual’s probe video, it has to be compared with a gallery video of the subject he or she claims to be, to determine whether they depict the same person. Within each verification algorithm, a proximity/distance measure is defined to quantify the similarity between two videos. The claims about identity are rejected or accepted according to a threshold value learned for said similarity metrics. We incur in a false positive if a false claim is accepted and a false negative whenever a true claim is rejected. The false accept rate (FAR) is the percentage of false positives and the false reject rate (FRR) is the percentage of false negatives encountered while running an algorithm. Both FAR and FRR provide a measure of the quality of an algorithm. A typical performance metrics for a recognition system is its receiver operating characteristic (ROC), which plots the FAR against the FRR as the similarity threshold varies. More synthetically, the verification rate is often employed, defined as the ratio of correct decisions (true positives and true negatives) divided by the total number of decisions. This metrics is more synthetic, however it is only significant if the number of false and true claims to be verified are the same or at least comparable. An over-trusting verification system, programmed to always output a yes decision, in a world of completely honest people, would achieve a 100 % verification rate, although it should be noted that such a world doesn’t exist and such a system is clearly useless.

**Closed-set identification** requires us to identify the person in a probe video by comparing it against gallery videos of different people and decide for the most similar one. We are guaranteed *a priori* that the subject in the probe video is depicted in at least one gallery video. Performance is commonly measured through the correct identification rate, counting the number of correct identifications divided by their total amount. Unlike verification, identification becomes inherently more difficult as the number of people in the gallery increases.



**Open-set identification** differs from closed-set identification in that we are not guaranteed that the probe subject appears in the gallery, making it a harder task. The correct identification rate can once more be defined as the number of correct decisions divided by the number of total decisions. This metric may vary significantly for the same algorithm depending on the “openness” of the dataset, which relates to the percentage of probe videos that depict a subject external to the gallery.

**Clustering** is a problem defined on a collection of videos, which have to be grouped based on the people they depict. The number of people in the videos can be undefined or known *a priori*, the latter case being easier to decide. The correct outcome of a clustering procedure is a partition of the video collection by subject.

In this document we test our novel dynamical model on tasks of verification, closed-set and open-set identification, showing the versatility it has inherited from PLDA while systematically outperforming it.

## 1.2 Peculiarities and challenges

Inference about identification from still images and from videos share the same problematic factors such as:

- Pose variations: the same subject can be captured from different angles, hence correspondence between pixels and face features is not guaranteed across pictures.
- Illumination variations: pictures can be lit by sources of different magnitude and orientation.
- Expression variations: over 30 facial muscles are responsible for countless configurations (expressions), which alter the basic appearance of a human face.
- Occlusions: when viewing conditions are not controlled, obstacles can come between the subject and the sensor, covering the former during acquisition.
- Motion blur: depending on the exposure time and the relative speed of the subject and the sensor, face pictures may be fuzzy and hard to identify even for humans.

In tackling these nuisances, video-based algorithms can exploit some peculiarities of the video realm. For example, multiple frames are available and the information carried by each of them can be combined in different manners. Depending on the application and its properties, robustness to illumination variations or improved performance through 3D modeling can be achieved taking into account the frame set as a whole. Furthermore, video sequences contain temporal information about the proximity of frames in time and consecutive frames are expected to be more correlated than frames that are far apart in time. Research in face recognition from videos is often based on research about still image identification, exploiting the peculiar advantages of video sequences to increase performance. This thesis is in fact an example of said approach, factoring temporal dependencies into PLDA models, which will be introduced in the next chapter.

## 1.3 Approaches

The authors of [3] provide an interesting taxonomy for face recognition in videos which we borrow and synthesize in figure 1.1. The classification is determined by how the additional information from multiple frames is exploited. Two essentially different approaches are

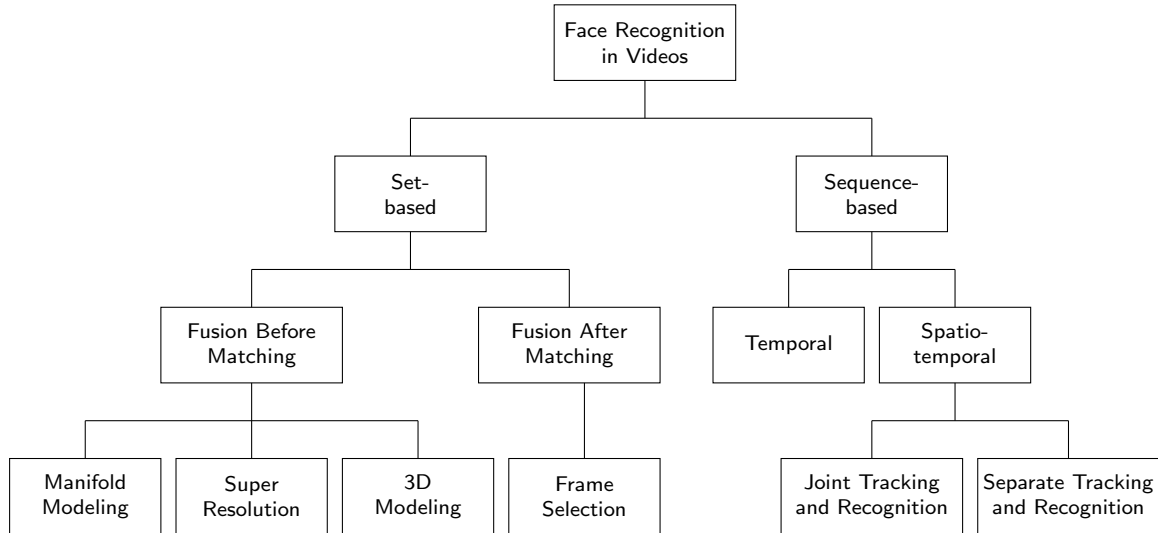


Figure 1.1: Taxonomy for face inference in videos. This figure shows when and how the information from multiple frames is exploited within a recognition pipeline. The procedure proposed in this thesis can be ascribed to Sequence-based, Spatio-temporal, Separate Tracking and Recognition.

set-based and sequence-based. **Set-based** techniques treat a video as a bag-of-images, completely ignoring temporal information. Sequence-based algorithms, on the other hand, exploit the ordering of the frames and their mutual position within a video to draw conclusions about the identity of the depicted subject.

Within the set-based approaches the information from different frames can be fused **before matching** or **after matching**. In the latter case, the matching procedure exploited is essentially one of face-recognition in images, generalized to videos by fusing the image-level results into the video-level domain, for instance through majority voting.

**Frame selection** allows us to pick a subset of all the video frames which is deemed sufficient to characterize the appearance of a subject. A key idea is that low quality frames and redundant frames can be discarded, improving recognition while compressing the available information. Also, gallery frames could be selected in order to span a wide range of illumination and pose conditions in view of uncontrolled and varying conditions in probe frames to be matched.

In [5], Berrani and Garcia employ robust PCA to select the frames that maximize a robust measure of variance within a low-dimensional base of eigenvectors onto which pixel values are projected. This criterion rules out frames that were captured under peculiar pose/illumination/expression, which lie too far from the others in the new space and are hence considered outliers.

Frame selection generalizes to frame weighting, wherein the importance of each frame in the decision process is determined by a continuous value. In [38] Zhang and Martinez propose a weighted selection scheme that exploits gaussian mixtures to find similarities within video frames, while assigning low weights to low-quality images that might be disruptive for the matching process.

Stallkamp et al. [27] devised a selection method that takes into account two different distance metrics and their combination for frame selection. The distance-to-model metric excludes probe frames whose distance from their closets match in the gallery is larger than a certain threshold, possibly due to its peculiar illumination and pose conditions.

The distance-to-second closest is a measure devised to systematically leave out ambiguous frames whose best and second-best matches are tied, thus extrinsically harder to decide between.

Another key idea for frame selection is to model and encourage intra-class diversity within the gallery, guaranteeing that the available frames within will be representative of an individual under most illumination/pose/expression configurations.

Clustering [13, 33, 34] is often exploited to group frames that share similar configurations. Matching can then be carried out as a subset-to-subset comparison or by selecting one significant exemplar from each cluster, for instance the cluster mean, which is regarded as significant of the whole subset. Exemplar-based approaches only require a limited amount of comparisons yielding a significant computational advantage while retaining most useful information.

In [13] Hadid and Pietikäinen employ a locality preserving LLE to reduce dimensionality while preserving similarity, prior to applying clustering. The cluster centers are regarded as a significant base for intra-class variations and taken into account for robust matching.

Another interesting approach consists of specifically encouraging the selection of high quality frames. Thomas et al. [28] introduce a “faceness” measure produced by the recognition software *FaceIt*. This measure is exploited for different strategies, by ordering video frames according to their faceness and selecting for instance the first  $N$  frames, or  $N$  evenly spaced frames from the first  $M$  with highest faceness.

In [36] Xiong and Jaines introduce an algorithm to systematically select the most suitable frames, based on their intrinsic and extrinsic properties. The intrinsic quality of a face is determined by its orientation, scale and illumination. The extrinsic quality is related to clustering so that a new frame is accepted if, in relation with the other frames previously picked, it achieves limited intra-class variance and increased inter-class distance.

**Super-resolution** methods tackle the problem of low image quality by incorporating information from multiple frames to reconstruct the high-frequency content lost during acquisition in challenging settings. The problem is typically formulated by modelling an image degradation process that has to be inverted in order to restore the high quality image:

$$\mathbf{y}(t) = \mathbf{M}(t)\mathbf{x}(t) + \mathbf{w}(t). \quad (1.1)$$

Here  $\mathbf{y}(t)$  denotes the observed image sequence up to a time step  $t$ ,  $\mathbf{M}(t)$  models the corruption process due to motion blur and sensor limitation,  $\mathbf{x}(t)$  represents the original, high quality sequence of frames to reconstruct, and  $\mathbf{w}(t)$  accounts for the remaining sources of noise. The reconstruction is typically carried out by minimizing a cost function consisting of the square reconstruction error plus a regularization term enforcing smoothness and mitigating the influence of noise.

The iterative back-projection algorithm, originally applied to face recognition by Zhou and Bhanu [39], is a classical super-resolution technique, whereby edges are recovered and enhanced by subtracting from the original images a blurred version of themselves. Al-Azzeh et al. [1] applied this procedure in combination with an efficient frequency-based alignment procedure, and completed their recognition pipeline with PCA to test their algorithm in a more scalable setting.

The application of PCA to the recovered images offers no guarantee that the high quality previously recovered will be retained. Furthermore the inversion of matrix  $\mathbf{M}(t)$  poses a significant computational burden in the pixel domain, whereas performing it in the PCA subspace results in lower computational requirements. For this reason, Gunturk et al. [12] transfer the problem from the pixel domain to the subspace spanned by the principal

---

components. In the low-dimensional domain, equation (1.1) is further complicated by an error term accounting for PCA reconstruction error.

It should be noted that the taxonomy in figure 1.1 is introduced for exposition convenience and the recognition schemes listed therein are not mutually exclusive. For instance, Jilela and Ross [15] combine super-resolution with frame selection by ruling out the frames corrupted by motion blur. In order to do so, they devise an intra-frame motion parameter  $\beta$  which measures intensity differences between points that are aligned in two consecutive frames. If  $\beta$  is smaller than a certain threshold, no significant movement is assumed to occur between two consecutive frames, which are hence selected for the reconstruction of a high quality image, otherwise they are discarded.

**3D Modeling** specifically copes with the problem of varying pose and correspondences between pixels and features of the subjects' face. 3D models are often exploited to generate synthetic images of gallery and probe subjects under matching pose and illumination conditions. Park et al. [22], for example, build offline a 3D model for each subject in the gallery. They subsequently exploit an SVM to estimate pose and illumination for the subject in the probe, and accordingly rotate and illuminate each gallery subject in an artificial fashion. Finally a 3D-to-2D projection is carried out, effectively synthesizing videos for each gallery subject under pose and illumination matching the ones of the probe. Alternatively, a 3D model can be learnt from the probe video, which is then rendered under varying modelled conditions corresponding to videos from the gallery.

An alternative to synthesis-driven schemes are model-comparison techniques, whereby 3D models are learnt from both gallery and probe and directly compared. An example of this approach is provided by Liu and Chen [19], who model the human head as a 3D ellipsoid, upon which people's faces are back projected, resulting in a texture map. Said map is divided into deformable patches, each of which corresponds to a specific region of the head. PCA is applied within each patch separately. The distance between 2 face models is then given by the sum of the distances between each pair of patches and exploited for recognition decisions.

**Manifold modelling** techniques characterize the possible configurations of a human face without explicitly modelling generative processes. The manifold associated to person  $i$  can be expressed as

$$\mathbf{X}_i = \{c(\mathbf{x}_i) \mid c \in C\}, \quad (1.2)$$

where  $C$  is the set of all possible configurations and  $c(\mathbf{x}_i)$  denotes the appearance of individual  $i$  under configuration  $c$ . The union of all  $\mathbf{X}_i$  yields the manifold of all faces.

Pioneering work was carried out in [37], where Yamaguchi et al. use PCA to represent faces of an individual as a linear subspace. The cosine of the principal angle between subspaces serves as a measure of similarity between faces, used in the decision process.

Other attempts were made in this direction, although linear subspaces struggle to characterize some inherent non-linearities such as illumination variations caused by light from multiple sources.

In [32] Wang et al. model face manifolds as a collection of linear patches, learnt with a method called Maximum Linear Patches. The number of patches is not decided in advance, rather it is determined adaptively. This technique iteratively constructs MLPs until the linearity constraint is broken. The mutual subspace angle is then used to quantify the distance between patches, which are then fused into a global distance metric for two faces.

Synthetic data can be used for manifold learning, as shown by Lina et al. in [18]. A global linear subspace is built through PCA. The available data is clustered according to pose angles and view-specific means and covariance matrices are learnt. First order and second order moments, super-imposed to the global linear subspace, characterize a

view-specific manifold. Artificial transformations such as blurring and slight rotations are used to learn more significant values for mean and covariance. The statistics for untrained poses are learnt through interpolation. Inference is performed based on the Mahalanobis distance between the learnt manifolds.

An interesting approach adopted by Arendjelović and Cipolla [2] is to consider face images as i.i.d. samples drawn from a low-dimensional PDF in the face space, transported by a mapping function onto a manifold in the image space. The authors resort to kernel PCA to unfold the face manifold, allowing them to model non-linearities while retaining computational efficiency. Finally a symmetrized version of the Kullback-Leibner distance (the Resistor-Average distance) is employed to perform inference about identity.

So far we have concentrated on set-based approaches, which combine the information from multiple frames disregarding their order. **Sequence-based** approaches, on the other hand, leverage the temporal information naturally available with videos. Dynamics is useful for face tracking purposes as well as for learning idiosyncratic variations peculiar to each individual. Note that both the notions of tracking and variation require a well defined time ordering of the frames in a video.

**Spatio-temporal** schemes combine static spatial information with dynamical cues in order to increase the performance of a recognition system. Algorithms ascribed to this category can be further divided, according to whether they perform joint face tracking and recognition or whether they split the two tasks.

**Separate tracking and recognition** offers the advantage of flexibility in that the tracking and the recognition algorithm can be devised independently of each other. The algorithm proposed in this thesis purely concentrates on the problem of face recognition, exploiting both temporal and spatial information in the process. The algorithm itself requires pre-aligned faces as an input, hence we integrate it into a recognition pipeline, wherein it is preceded by a step of independent face tracking and alignment. For this reason our recognition pipeline can be ascribed to this category of procedures.

A common technique to learn the statistics associated with an individual's face while comprising temporal information are Hidden Markov Models. Based on video frames, which are regarded as the output of an HMM  $\Lambda_i = (\mathbf{A}_i, \mathbf{B}_i, \boldsymbol{\pi}_i)$ , the model parameters are learnt through the Viterbi algorithm. A different model  $\Lambda_i$  is associated with each person appearing in the gallery. When a probe sequence  $\mathcal{O}$  is available, it is ascribed to the identity  $i$  whose HMM  $\Lambda_i$  maximizes the probability of the sequence conditional on the learnt parameters  $(\mathbf{A}_i, \mathbf{B}_i, \boldsymbol{\pi}_i)$ . A significant example of this approach was devised by Liu and Cheng [20], who designed an adaptive HMM capable of online learning.

In [11] Gorodnichy resorts to an autoassociative neural network that accumulates information over time by updating its synapses. Based on physiological considerations, the author emphasizes the dynamical nature of human face recognition. Synaptic plasticity helps us to memorize from sequences of subsequent stimuli, rather than from isolated face frames. Analogously the author proposes a scheme to update a neural network's synaptic weights in the presence of frame sequences ordered in time.

Aggarwal et al. [10] propose the first recognition approach based on an autonomous ARMA model. Within their framework, observations correspond to video frames and hidden states are essentially related to pose. Analogously to [20], a separate model is learnt for every individual in the gallery. Finally, principal angles between learnt models are employed as a discrimination metrics.

Feature extraction is commonly carried out at a frame level, exploiting the temporal information in the subsequent steps. Hadid and Pietikäinen [14] introduce the Volume Local Binary Patterns (VLBP), to extract features from 3D spatio-temporal neighbourhoods

---

comprising neighbouring pixels from neighbouring frames. Recognition is performed by comparing feature histograms of VLBP and subsequently exploiting an SVM classifier.

**Joint tracking and recognition** permits a bidirectional flow of information between tracker and recognition system, which can increase the overall performance if properly leveraged. The tracker is likely to return faces that comply with the appearance model employed by the recognition system. Also, information about the pose can be shared seamlessly between the two components of the identification pipeline.

Li et al. [17] propose a unified tracking and recognition scheme which models face appearance as a probabilistic manifold which is approximated by a collection of linear subspaces, each of which is associated with a different pose. The pose information is used by the recognition system to perform identity inference at a frame level. The frame-specific information about identity is fused at a video level with a bayesian approach. At the same time the same pose information from a frame is used for tracking in the next one, wherein the pose is modelled as a gaussian distribution supposed to be centered at its previous value. Connectivity of pose-dependent subspaces is thus modelled with transition matrices that enforce proximity. This scheme is interesting as recognition and tracking share the same appearance model, reducing the misalignments between tracker output and recognition input.

**Temporal methods** concentrate on the peculiar variational patterns of a face rather than its static aspect or mean appearance. Temporal face cues cannot be disguised as easily as appearance since they are part of nearly subconscious mechanisms. At the same time, they have been shown to be fundamental for human perception in recognizing familiar faces and, more in general, in complementing static information in the presence of low quality video sequences.

Chen et al. [7] propose a temporal method based on motion flow fields feature extraction. The motion flow field is a classical dynamical feature related to the temporal evolution of a depicted object. Videos of subjects uttering the same word are acquired in a constrained setting, properly synchronized and landmarked. An optical flow is computed for every frame and concatenated in a high-dimensional representation of an individual. PCA and LDA are exploited as a dimension reduction tool yielding individual representations which are compared through their euclidean distance.

Matta and Dugelay [21] propose a system based on piece-wise head motion measurements. The displacement of an individual's eyes, nose and mouth over frame sequences is computed and normalized for robustness with respect to rotation and scaling. These features are employed to train individual GMMs which are used as a bayesian classifier.

Benedikt et al [4] concentrate on the temporal data from videos of people uttering the word "Puppy". Faces from different videos are aligned and cropped around the region of the lips. Further dimensionality reduction is achieved through PCA, and Weighted Dynamic Time Warping (WDTW) is employed as a pattern matching measure.

Temporal schemes have shown promising performance and good robustness to potentially deceptive face ornaments, however they remain a vastly unexplored domain less mature than spatio-temporal techniques. In this document we present a novel generative model (DPLDA) that can be ascribed to this latter paradigm.

- In chapter 2 we present a review of PLDA, which constitutes the static foundation of DPLDA, originally conceived for recognition in still images.
- The first part of chapter 3 presents a dynamical generalization of PLDA. The phases of learning and recognition are then described in detail, with the respective implementation details.

- In chapter 4 we introduce the remaining components of our recognition pipeline, and present the results it achieves on different databases for different recognition tasks.
- Conclusions and further directions of research are outlined in chapter 5.





## 2. PLDA

### 2.1 A preliminary review of LDA

Linear discriminant analysis (LDA) [9] is the foundation for the generative models described throughout this document. LDA is a classical technique employed in pattern recognition to separate a set of observations into two or more classes. Similarly to PCA, LDA is a deterministic dimensionality reduction tool that specifically aims at retaining most useful information (i.e. the variance) in the original dataset. However, unlike PCA, it is specifically designed to take into account the different classes that training samples are ascribed to, and make sure the data is projected onto some linear subspace that guarantees good discriminability between the different classes.

Intuitively, we would like samples of the same class to look very similar - short distance in the projection subspace - and samples from different classes to be far apart - long distance in the projection subspace. This idea can be formalized by defining the within-class scatter matrix for  $C$  different classes

$$\mathbf{S}_W = \sum_{j=1}^C \frac{1}{N_{c_j}} \sum_{\mathbf{x}_i \in c_j} (\mathbf{x}_i - \boldsymbol{\mu}(c_j)) (\mathbf{x}_i - \boldsymbol{\mu}(c_j))^T, \quad (2.1)$$

where  $\mathbf{x}_i$  are the observations and  $\boldsymbol{\mu}(c_j)$ ,  $j \in (1, C)$  are the centers for each class.  $\mathbf{S}_W$  captures the variance for data with the same label, whereas the between-class scatter matrix, defined as

$$\mathbf{S}_b = \sum_{j=1}^C (\boldsymbol{\mu}(c_j) - \boldsymbol{\mu}) (\boldsymbol{\mu}(c_j) - \boldsymbol{\mu})^T, \quad (2.2)$$

measures the overall distance between classes. We then look for the projection matrix  $\mathbf{W}$  that solves the optimization problem

$$\max \text{Tr}[\mathbf{W}^T \mathbf{S}_b \mathbf{W}] \text{ s.t. } \mathbf{W}^T \mathbf{S}_W \mathbf{W} = 1, \quad (2.3)$$

which yields a new representation for the data,  $\mathbf{y}_i = \mathbf{W} \mathbf{x}_i$ , featuring both good discrimination power and low dimensionality.

### 2.2 PLDA: a static model for inference about identity

Probabilistic LDA (PLDA), as detailed in [16], represents a probabilistic, generative framework for inference about identity in still images. Let us assume we have a training set consisting of  $J$  images per individual, with a total of  $I$  different individuals. Then, according to PLDA, the generative model behind the  $j$ -th image of the  $i$ -th individual is

$$\mathbf{x}_{ij} = \boldsymbol{\mu} + \mathbf{F} \mathbf{h}_i + \mathbf{G} \mathbf{w}_{ij} + \boldsymbol{\epsilon}_{ij}. \quad (2.4)$$

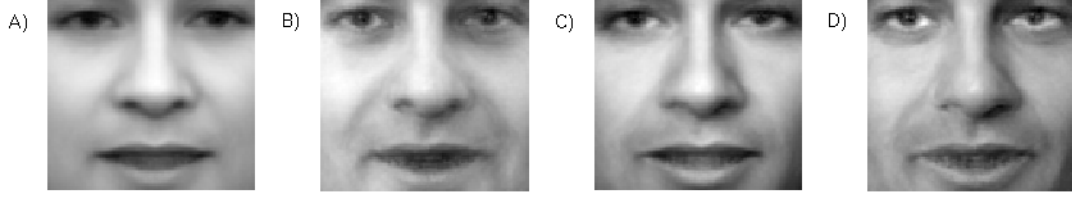


Figure 2.1: Generative process. A final image D) can be decomposed into a generic mean face  $\mu$  (A), a component from the identity subspace  $\mu + \mathbf{F}\mathbf{h}_i$  (B) and a component from the pose/illumination/expression subspace  $\mu + \mathbf{G}\mathbf{w}_{ij}$  (C).

Here  $\mathbf{x}_{ij}$  is the image itself, represented by a vector of pixel values. According to the generative model, it depends on  $\mathbf{h}_i$ , which is the latent variable specific to the  $i$ -th person, and on  $\mathbf{w}_{ij}$  which is the latent variable associated with the setting in which the picture was taken. Both  $\mathbf{h}_i$  and  $\mathbf{w}_{ij}$  are unobserved random variables. Intuitively,  $\mathbf{h}_i$  captures the face features of a person that consistently determine his/her look, whereas  $\mathbf{w}_{ij}$  represents incidental conditions such as pose, illumination and expression that influenced the picture at the moment it was taken.

More explicitly, we can look at the model from a generative point of view, according to which an image is created by taking a realization of random variables  $\mathbf{h}_i$  and  $\mathbf{w}_{ij}$ , projecting them according to matrices  $\mathbf{F}$  and  $\mathbf{G}$ , adding the mean vector and a sample from the noise process. The adoption of said model is not motivated by a belief that pictures of individuals are generated in strict compliance to this equation, rather by a combination of simplicity and good discriminatory power. In figure 2.1 we depict a breakdown of a face image into its generative components. Note that the shade in the bottom-right corner is correctly ascribed to  $\mathbf{G}\mathbf{w}_{ij}$  as an incidental illumination condition.

In more rigorous terms,  $\mathbf{F}$  is a factor matrix whose columns span the between-individual subspace. Each individual is assumed to have a unique position in said subspace, which sets it apart from everyone else and is represented by the hidden variable  $\mathbf{h}_i$ . Analogously,  $\mathbf{G}$  is a matrix whose columns span the within-individual subspace.  $\mathbf{w}_{ij}$  represents the position of image  $\mathbf{x}_{ij}$  therein and is responsible for the differences in photos of the same individual. Vector  $\mu$  is simply the mean of all images, which is computed initially for detrending purposes and  $\epsilon_{ij}$  is stochastic noise that justifies the observations' distance from the hypothesized model. Said noise is supposed to be a random variable with zero mean and (anisotropic) diagonal covariance matrix. This feature is very important for the discriminative power of the model, allowing it to learn which pixels carry more significant information for the purposes of identity inference. Figure 2.2 portrays different positions in the above-mentioned subspaces, together with the mean face  $\mu$  and a sample of  $\epsilon_{ij}$ .

The prior and conditional distributions associated to model (2.4) are defined as follows:

$$Pr(\mathbf{x}_{ij}|\mathbf{h}_i, \mathbf{w}_{ij}, \theta) = \mathcal{N}(\mathbf{x}_{ij}|\mu + \mathbf{F}\mathbf{h}_i + \mathbf{G}\mathbf{w}_{ij}, \Sigma), \quad (2.5)$$

$$Pr(\mathbf{h}_i) = \mathcal{N}(\mathbf{h}_i|\mathbf{0}, \mathbf{I}), \quad (2.6)$$

$$Pr(\mathbf{w}_{ij}) = \mathcal{N}(\mathbf{w}_{ij}|\mathbf{0}, \mathbf{I}), \quad (2.7)$$

where  $\mathcal{N}(\mathbf{x}|\mathbf{a}, \mathbf{B})$  indicates a gaussian PDF on random vector  $\mathbf{x}$  with mean  $\mathbf{a}$  and variance matrix  $\mathbf{B}$ . Furthermore the set of all the parameters governing the models is defined as  $\theta = (\mu, \mathbf{F}, \mathbf{G}, \Sigma)$ . The adoption of gaussian PDFs, in conjunction with the linear equation



Figure 2.2: Different directions in both subspaces. Some directions in the between-individual subspace (B), capturing the features of specific individuals. Some directions in the within-individual subspace (D), capturing differences in expression and illumination. Overall mean face (A). A sample from noise process  $\epsilon_{ij}$  (C)

(2.4), determine nice closed-form solutions for Bayes law computations, which are at the core of generative models.

The term  $\mu + \mathbf{F}\mathbf{h}_i$  accounts for between-individual variance, whereas  $\mathbf{G}\mathbf{w}_{ij}$  is responsible for within-individual variance. From this perspective it is possible to understand the claim that “the relationship between PLDA and standard LDA is similar to that between factor analysis and PCA”, formulated by the authors of [16] without any further explanation.

PCA poses the problem of finding a projection matrix from a high-dimensional space to a smaller one, that, given a series of high-dimensional observations, allows us to compress them while retaining the most possible variance. Factor analysis, conversely, is based on a probabilistic framework featuring a low-dimensional (unobserved) latent variable, assumed to be gaussian with zero mean, in analogy with PLDA. Given a set of observations, the procedure aims at finding the projection matrix from the low-dimensional space of the hidden variables to the high-dimensional space of the observations, that maximizes the likelihood of the observed variables under this model. It has been shown in [29] that, as the variance of the latent variable in factor analysis tends to zero, the optimal projection matrix converges to the transpose of the PCA solution, thus giving a precise relationship between the two frameworks. PCA is extremely effective for tasks of unsupervised learning, but it fails to take into account the labels (identities) that come with our training set.

When the goal is to classify data into different classes, LDA is a handy dimensionality reduction tool. In particular, LDA maximizes between-class variance while minimizing within-class variance by finding a suitable projection matrix from the high-dimensional observation space onto a smaller one. In other words, data from the same class is clustered together and kept distant from other classes, thus favouring discrimination. The concepts of between-class and within-class variance entail a notion of labels (classes) that LDA naturally takes into account. PLDA, on the other hand, resembles factor analysis in that it is based on a mapping from a low-dimensional latent space to the high-dimensional data space. Training a PLDA model can be seen as finding the 2 factor matrices  $\mathbf{F}$  and

---

$\mathbf{G}$  that fit the observations optimally as a projection from the smaller subspace to which the couple  $(\mathbf{h}_i, \mathbf{w}_{ij})$  belongs. Furthermore, just like LDA, PLDA clusters together images  $\mathbf{x}_{\bar{i}j}$  from the same class  $\bar{i}$  by taking into account the respective latent variables  $(\mathbf{h}_{\bar{i}}, \mathbf{w}_{\bar{i}j})$ . This clustering is achieved through the hard constraint that  $\mathbf{h}_{\bar{i}}$  should be the same across all pictures of person  $\bar{i}$ , thus restricting images of the same individual to vary less in the subspace spanned by the columns of  $\mathbf{F}$  and  $\mathbf{G}$ . Regardless of this loose explanation, to the best of our knowledge, no precise mathematical characterization relates LDA with PLDA. We end this section with a list of four assumptions that enclose the essence of PLDA:

- **Face images depend on several interacting factors**, such as identity, pose, expression, illumination.
- **Image generation is noisy**, due to both sensor noise and model imprecision.
- **Identity cannot be known exactly**: embracing a bayesian point of view, no attempt will be made to estimate identity variables. However:
- **Recognition tasks do not require identity estimates**: we can decide whether two images share the same identity without explicitly estimating said identity.

The discussion carried out so far has expanded on the two top assumptions, whereas to better understand the two bottom ones we need to delve deeper into the recognition technique associated to PLDA, which is the subject of the next section.

## 2.3 Recognition

Suppose, for the moment, that a set of parameters  $\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{F}, \mathbf{G}, \boldsymbol{\Sigma})$  has been learnt according to some optimality criterion. Recognition tasks typically require us to determine which images depict the same person, by exploiting our knowledge about how the pictures were created. Within our framework, recognition implies estimating which images  $\mathbf{x}_1 \dots \mathbf{x}_k$  were generated by the same identity latent variable  $\mathbf{h}$  in a model governed by parameter set  $\boldsymbol{\theta}$ . To avoid notational ambiguity, let us define the *process*  $\mathcal{P}$  behind a set of images  $(\mathbf{x}_1 \dots \mathbf{x}_N)$  as the ground truth of the people  $(\mathbf{h}_1 \dots \mathbf{h}_N)$  depicted in those pictures. When presented with a task of recognition, or inference about identity, we are given a set of images, with partial knowledge of the process that generated them (the labels of the training set) and we're asked to determine the missing information about said process (the labels of the test set). Note the distinction between our definitions of process and model. The process actually contributes to the generation of the data set, the model is an artificial description we adopt to characterize how this contribution takes place.

Since PLDA is a probabilistic generative model, a natural criterion for recognition is given by Maximum a posteriori (MAP) estimation. Given a set of images  $\mathbf{x}_{1\dots N}$ , we want to find the process  $\mathcal{P}^*$  that best explains the observation:

$$\mathcal{P}^* = \arg \max_{\mathcal{P}} Pr(\mathcal{P} | \mathbf{x}_{1\dots N}), \quad (2.8)$$

where

$$Pr(\mathcal{P} | \mathbf{x}_{1\dots N}) \propto Pr(\mathbf{x}_{1\dots N} | \mathcal{P}) P(\mathcal{P}). \quad (2.9)$$

In this regard, the effort lays in quantifying the likelihoods  $Pr(\mathbf{x}_{1\dots N} | \mathcal{P})$ , whereas, for the process priors  $P(\mathcal{P})$ , uniform probabilities are assumed.

Recognition can be performed for slightly different tasks, among which closed set face identification, open set identification, clustering and verification, as described in section

1.1. One of the advantages of PLDA is that it allows us to tackle all these problems within the same framework. Let us consider, for example, an open set face identification problem, where  $N$  images  $\mathbf{x}_1, \dots, \mathbf{x}_N$  of  $N$  different people represent the gallery and we need to determine whether a probe image  $\mathbf{x}_p$  shares its identity variable  $\mathbf{h}$  with one of them. In order to do this, we will compute and compare  $N + 1$  different likelihoods, conditional on all the possible processes that may have generated the observation.

$$Pr(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}_p | \mathcal{P}_0) = \prod_{n=1}^N Pr(\mathbf{x}_n) Pr(\mathbf{x}_p), \quad (2.10)$$

$$Pr(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}_p | \mathcal{P}_i) = \prod_{n=1, n \neq i}^N Pr(\mathbf{x}_n) Pr(\mathbf{x}_p, \mathbf{x}_i), \quad i = 1 \dots N \quad (2.11)$$

Here  $\mathcal{P}_0$  denotes the process where the person in the probe image does not appear in the gallery whereas hypotheses  $\mathcal{P}_i$  assume that images  $(\mathbf{x}_p, \mathbf{x}_i)$  depict the same person. Terms  $Pr(\mathbf{x}_n)$  and  $Pr(\mathbf{x}_p)$  in equation (2.10) represent the likelihood of a single image and can be evaluated according to

$$Pr(\mathbf{x}_n) = \iint Pr(\mathbf{x}_n, \mathbf{h}_n, \mathbf{w}_n) d\mathbf{h}_n d\mathbf{w}_n, \quad (2.12)$$

$$Pr(\mathbf{x}_p) = \iint Pr(\mathbf{x}_p, \mathbf{h}_p, \mathbf{w}_p) d\mathbf{h}_p d\mathbf{w}_p. \quad (2.13)$$

The computation is carried out by comprising latent variables into a joint probability and then marginalizing them out with a bayesian approach. Note this is in perfect agreement with the third and fourth assumptions listed at the end of the previous section. In fact we estimate a PDF for  $\mathbf{h}$  but we're not interested in the single most plausible value.

$Pr(\mathbf{x}_p, \mathbf{x}_i)$  represents the likelihood of images  $\mathbf{x}_p$  and  $\mathbf{x}_i$  under the hypothesis that they share the same identity hidden variable:

$$Pr(\mathbf{x}_p, \mathbf{x}_i) = \iiint Pr(\mathbf{x}_p, \mathbf{x}_i, \mathbf{h}_i, \mathbf{w}_p, \mathbf{w}_i) d\mathbf{h}_i d\mathbf{w}_p d\mathbf{w}_i. \quad (2.14)$$

The constraint of sharing the same identity is expressed by integrating just along  $\mathbf{h}_i$ , implicitly assuming  $\mathbf{h}_p$  to have the same value. Note, once more, that we integrate along the whole range of possible values of  $\mathbf{h}_i$ , without performing a specific estimate.

Equations (2.12), (2.13), (2.14) can be computed by solving the more general problem of evaluating the likelihood of  $N$  images  $\mathbf{x}_1, \dots, \mathbf{x}_N$  sharing the same identity variable  $\mathbf{h}$ , as explained below. Combining  $N$  generative equations we have

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_N \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu} \\ \vdots \\ \boldsymbol{\mu} \end{bmatrix} + \begin{bmatrix} \mathbf{F} & \mathbf{G} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{G} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{F} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{G} \end{bmatrix} \begin{bmatrix} \mathbf{h} \\ \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_N \end{bmatrix} + \begin{bmatrix} \boldsymbol{\epsilon}_1 \\ \boldsymbol{\epsilon}_2 \\ \vdots \\ \boldsymbol{\epsilon}_N \end{bmatrix}, \quad (2.15)$$

which can be written synthetically as

$$\mathbf{x}' = \boldsymbol{\mu}' + \mathbf{D}\mathbf{y} + \boldsymbol{\epsilon}'. \quad (2.16)$$

The random vectors  $\mathbf{x}'$  and  $\mathbf{y}$  are distributed according to the following prior and conditional PDFs:

$$Pr(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \boldsymbol{\theta}, \mathbf{I}), \quad (2.17)$$

---


$$Pr(\mathbf{x}'|\mathbf{y}) = \mathcal{N}(\mathbf{x}'|\boldsymbol{\mu}' + \mathbf{D}\mathbf{y}, \boldsymbol{\Sigma}'), \quad (2.18)$$

where

$$\boldsymbol{\Sigma}' = \begin{bmatrix} \boldsymbol{\Sigma} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \mathbf{0} & \cdots & \boldsymbol{\Sigma} \end{bmatrix}. \quad (2.19)$$

Note that model (2.16) has the form of a factor analyzer. From the properties of conditional gaussian distributions it's easy to show that

$$Pr(\mathbf{x}_{1...N}) = Pr(\mathbf{x}') = \mathcal{N}(\mathbf{x}'|\boldsymbol{\mu}', \boldsymbol{\Sigma}' + \mathbf{D}\mathbf{D}^T). \quad (2.20)$$

Having identified this probability distribution, it's now easy to compare a probe image against all the gallery images, evaluate (2.12), (2.13), (2.14) which in turn allow us to compute (2.10) and (2.11) and combining them with the uniform priors over the processes, we can select the most plausible process by picking the option that maximizes the posterior.

## 2.4 Learning

We postponed to this section a treatment of how the model parameters  $\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{F}, \mathbf{G}, \boldsymbol{\Sigma})$  can be evaluated in order to perform recognition. If we knew the hidden variables  $\mathbf{h}_i$  and  $\mathbf{w}_{ij}$ , we could employ a maximum-likelihood framework to estimate a sensible value for  $\boldsymbol{\theta}$ . Likewise, if  $\boldsymbol{\theta}$  was known, evaluating the hidden variables would be easy. When presented with a problem like this one, seemingly recursive and unsolvable, a typical solution is resorting to the Expectation-Maximization (EM) algorithm. This procedure allows us to find maximum likelihood solutions in models comprising latent variables, hence it perfectly applies to our case. A formal characterization of the EM algorithm and its properties can be found in appendix A. Here we limit our scope to an intuitive explanation of the steps and the reasoning behind this procedure.

Let us indicate with  $\mathbf{X}$  the totality of the observations and with  $\mathbf{Z}$  the underlying latent variables. Our goal is to find the parameter  $\boldsymbol{\theta}$  that maximizes

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}), \quad (2.21)$$

or equivalently its logarithm. It often happens that direct optimization on  $p(\mathbf{X}|\boldsymbol{\theta})$  is a hard problem, whereas optimization on the joint distribution  $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$  is considerably easier. To make this claim more tangible, we can relate it to PLDA: within our current framework, direct computation of  $p(\mathbf{X}|\boldsymbol{\theta})$  implies evaluation of the inverse for covariance matrix in equation (2.20), which takes the form

$$(\boldsymbol{\Sigma}' + \mathbf{D}\mathbf{D}^T)^{-1} = \begin{bmatrix} (\boldsymbol{\Sigma} + \mathbf{G}\mathbf{G}^T)^{-1} & & \\ & \ddots & \\ & & (\boldsymbol{\Sigma} + \mathbf{G}\mathbf{G}^T)^{-1} \end{bmatrix} + \begin{bmatrix} (\boldsymbol{\Sigma} + \mathbf{G}\mathbf{G}^T)^{-1}\mathbf{F} \\ \vdots \\ (\boldsymbol{\Sigma} + \mathbf{G}\mathbf{G}^T)^{-1}\mathbf{F} \end{bmatrix} (\mathbf{I} + \mathbf{J}\mathbf{F}^T(\boldsymbol{\Sigma} + \mathbf{G}\mathbf{G}^T)^{-1}\mathbf{F}) \begin{bmatrix} (\boldsymbol{\Sigma} + \mathbf{G}\mathbf{G}^T)^{-1}\mathbf{F} \\ \vdots \\ (\boldsymbol{\Sigma} + \mathbf{G}\mathbf{G}^T)^{-1}\mathbf{F} \end{bmatrix}^T, \quad (2.22)$$

followed by computation of the respective quadratic form and derivation with respect to  $\boldsymbol{\theta}$ . This is indicative of how hard direct optimization can be. On the other hand, evaluation

of  $\log[p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})]$  is usually simpler, since it can be broken down into the prior and conditional  $p(\mathbf{Z}|\boldsymbol{\theta})$ ,  $P(\mathbf{X}|\mathbf{Z}, \boldsymbol{\theta})$ , which can be separated additively thanks to the logarithm function. Also, here we see the first advantage of working with gaussian distributions, whose exponentials revert the logarithm, simplifying both out of the equation. To delve into the mathematical details for this example, see subsection 2.4.2.

Our efforts are hence devoted to maximizing the quantity  $\log[p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})]$  with respect to  $\boldsymbol{\theta}$ . A fundamental problem in this regard is that  $\mathbf{Z}$  is a latent variable and, as such, unknown. In the absence of certain, deterministic information about  $\mathbf{Z}$ , our next-best option is to evaluate its expectation and regard it as a reliable estimate for its real value. After this heuristic explanation, we can sketch the typical flow of the EM algorithm:

1. Initialize parameter  $\boldsymbol{\theta}^{old}$ .
2. **E step** - expectation: evaluate  $E(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})$
3. **M step** - maximization: update the value of  $\boldsymbol{\theta}$  by computing

$$\boldsymbol{\theta}_{new} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}), \quad (2.23)$$

where

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) = E_{\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}}[\log(p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}))] \quad (2.24)$$

is the expected value with respect to  $\mathbf{Z}$  of the joint log-probability defined in (2.21), conditioned on the available information  $(\mathbf{X}, \boldsymbol{\theta}^{old})$ .

4. Check for convergence

The alternation of E step and M step guarantees to increase the likelihood function  $p(\mathbf{X}|\boldsymbol{\theta})$  at every iteration (unless we reached a local maximum). A proof of this statement can be found in appendix A. In the next two subsections we provide a sketch of how the concepts we introduced apply to PLDA.

### 2.4.1 E step

Starting from a sequence of gallery images from the same person  $i$ , modeled according to equation

$$\begin{bmatrix} \mathbf{x}_{i1} \\ \mathbf{x}_{i2} \\ \vdots \\ \mathbf{x}_{iN} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu} \\ \vdots \\ \boldsymbol{\mu} \end{bmatrix} + \begin{bmatrix} \mathbf{F} & \mathbf{G} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{G} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{F} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{G} \end{bmatrix} \begin{bmatrix} \mathbf{h}_i \\ \mathbf{w}_{i1} \\ \vdots \\ \mathbf{w}_{iN} \end{bmatrix} + \begin{bmatrix} \boldsymbol{\epsilon}_{i1} \\ \boldsymbol{\epsilon}_{i2} \\ \vdots \\ \boldsymbol{\epsilon}_{iN} \end{bmatrix}, \quad (2.25)$$

we want to estimate the probability distribution of the respective latent variables, namely the within-individual variable  $\mathbf{w}_{ij}$ , specific to each image, and the identity variable  $\mathbf{h}_i$ , which is shared across the sequence. In a more compact notation, our aim is finding the sufficient statistics for  $\mathbf{z}_i$  from a model of the form

$$\mathbf{x}_i = \boldsymbol{\mu}' + \mathbf{D}\mathbf{z}_i + \boldsymbol{\epsilon}_i. \quad (2.26)$$

Exploiting the properties of gaussian distributions, it can be shown that

$$P(\mathbf{z}_i|\mathbf{x}_i) = \mathcal{N}(\mathbf{z}_i | (\mathbf{I} + \mathbf{D}^T \boldsymbol{\Sigma}'^{-1} \mathbf{D})^{-1} \mathbf{D}^T \boldsymbol{\Sigma}'^{-1} (\mathbf{x}_i - \boldsymbol{\mu}'), (\mathbf{I} + \mathbf{D}^T \boldsymbol{\Sigma}'^{-1} \mathbf{D})^{-1}), \quad (2.27)$$

with  $\Sigma'$  as defined in (2.19). As can be observed by looking ahead into equations (2.33) - (2.35), the sufficient statistics from this PDF are

$$E(\mathbf{z}_i|\mathbf{x}_i, \theta) = (\mathbf{I} + \mathbf{D}^T \Sigma'^{-1} \mathbf{D})^{-1} \mathbf{D}^T \Sigma'^{-1} (\mathbf{x}_i - \boldsymbol{\mu}'), \quad (2.28)$$

$$E(\mathbf{z}_i \mathbf{z}_i^T | \mathbf{x}_i, \theta) = (\mathbf{I} + \mathbf{D}^T \Sigma'^{-1} \mathbf{D})^{-1} + E(\mathbf{z}_i | \mathbf{x}_i, \theta) E(\mathbf{z}_i | \mathbf{x}_i, \theta)^T, \quad (2.29)$$

where we made explicit the dependency of the above values on the model parameters  $\theta$ . In order to compute the distribution of the latent variables  $\mathbf{z}_i$  associated to every person, the above calculations should be carried out for  $i = 1, \dots, I$ .

### 2.4.2 M step

In the M-step, we aim to update the model parameters  $\theta = (\boldsymbol{\mu}, \mathbf{F}, \mathbf{G}, \Sigma)$  in order to maximize the joint likelihood of observations and latent variables; however not knowing the exact value for the latent variables, we will maximize the expected value of the said joint likelihood, leveraging the results from the E-step. Consider the following model for a single image:

$$\mathbf{x}_{ij} = \boldsymbol{\mu} + [\mathbf{F} \quad \mathbf{G}] \begin{bmatrix} \mathbf{h}_i \\ \mathbf{w}_{ij} \end{bmatrix} + \epsilon_{ij}, \quad (2.30)$$

which can be rewritten as

$$\mathbf{x}_{ij} = \boldsymbol{\mu} + \mathbf{B} \mathbf{z}_{ij} + \epsilon_{ij}. \quad (2.31)$$

Our objective is finding

$$\theta^* = \arg \max_{\theta} Q(\theta, \theta^{old}), \quad (2.32)$$

where

$$\begin{aligned} Q(\theta, \theta^{old}) &= E_{\mathbf{Z}|\theta^{old}, \mathbf{X}} \{ \ln Pr(\mathbf{X}, \mathbf{Z}|\theta) \} \\ &= E_{\mathbf{Z}|\theta^{old}, \mathbf{X}} \left\{ \ln \prod_{i=1}^I \prod_{j=1}^J Pr(\mathbf{x}_{ij} | \mathbf{z}_{ij}) Pr(\mathbf{z}_{ij}) \right\} \\ &= -E_{\mathbf{Z}|\theta^{old}, \mathbf{X}} \left\{ \sum_{i=1}^I \sum_{j=1}^J \left[ \frac{1}{2} \ln |\Sigma| + \frac{1}{2} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{B} \mathbf{z}_{ij})^T \Sigma^{-1} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{B} \mathbf{z}_{ij}) \right. \right. \\ &\quad \left. \left. + \frac{D_h + D_w}{2} \ln 2\pi + \frac{1}{2} \mathbf{z}_{ij}^T \mathbf{z}_{ij} + \frac{f}{2} \ln 2\pi \right] \right\} \end{aligned}$$

Here we denote as  $D_w$  and  $D_h$  the dimensionality of the latent spaces and as  $f$  the size of the observations. Note, as anticipated, how the logarithm function in conjunction with gaussian PDFs render the computations simple to handle, allowing for nice closed form solutions. Indeed, by setting to zero the derivatives of  $Q$  with respect to  $\theta$ , we have:

$$\boldsymbol{\mu} = \frac{1}{IJ} \sum_{i=1}^I \sum_{j=1}^J \mathbf{x}_{ij} \quad (2.33)$$

$$\mathbf{B} = \left( \sum_{i=1}^I \sum_{j=1}^J (\mathbf{x}_{ij} - \boldsymbol{\mu}) E[\mathbf{z}_{ij}^T] \right) \left( \sum_{i=1}^I \sum_{j=1}^J E[\mathbf{z}_{ij} \mathbf{z}_{ij}^T] \right)^{-1} \quad (2.34)$$

$$\Sigma = \frac{1}{IJ} \sum_{i=1}^I \sum_{j=1}^J \text{diag} \{ (\mathbf{x}_{ij} - \boldsymbol{\mu})(\mathbf{x}_{ij} - \boldsymbol{\mu})^T - 2\mathbf{B} E[\mathbf{z}_{ij}] (\mathbf{x}_{ij} - \boldsymbol{\mu})^T + \mathbf{B} E[\mathbf{z}_{ij} \mathbf{z}_{ij}^T] \mathbf{B}^T \} \quad (2.35)$$



These values are then re-plugged into equations (2.28), (2.29) to update the sufficient statistics, which in turn will allow to update  $\boldsymbol{\theta}$ , carrying on with this cycle until an appropriate convergence criterion is met.



## 3. Dynamical PLDA

### 3.1 PLDA for videos

PLDA has been used for face recognition in still images for a wide variety of databases and different settings. When it comes to applications of inference from videos, it can still be employed as a set-to-set (probe video to gallery video) matching tool that compares all the frames of a probe video against all the frames of a gallery video and determines the overall most likely identity through maximum voting or some similar global metrics. The set-based paradigm allows us to mitigate the effect of noise and non-optimal viewing conditions by smoothing these anomalies across all the available frames. In [33] and [34] this idea is employed by a preliminary pose-dependent clustering and a subsequent matching of same-cluster frames. Different global metrics are explored and compared for the final matching scores.

This approach, despite achieving good results, ignores the temporal information that goes with a video, in particular the continuity that links the frames across time. It's reasonable to expect that two consecutive frames will look fairly similar; it's just as reasonable to impose proximity for the within-individual latent variables in two consecutive frames of the same video. Such property relies on a sequence-based approach, that specifically takes into account and models temporal dependencies and dynamics. We refer to this unprecedented approach as dynamical PLDA (or DPLDA). Throughout this chapter we provide the necessary mathematical framework and the implementation details to perform learning and inference in DPLDA models.

Consider the following formulation:

$$\mathbf{x}_{ij}^t = \boldsymbol{\mu} + \mathbf{F}\mathbf{h}_i + \mathbf{G}\mathbf{w}_{ij}^t + \boldsymbol{\epsilon}_{ij}^t, \quad (3.1)$$

$$\mathbf{w}_{ij}^t = \mathbf{A}_{ij}\mathbf{w}_{ij}^{t-1} + \mathbf{v}_{ij}^t, \quad (3.2)$$

$$\boldsymbol{\epsilon}_{ij}^t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}), \quad (3.3)$$

$$\mathbf{v}_{ij}^t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (3.4)$$

$$\mathbf{w}_{ij}^1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (3.5)$$

$$\mathbf{h}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (3.6)$$

The notation employed hereafter is explained in table 3.1. Note we didn't introduce it before to avoid confusion due to some overlapping with classical PLDA notation.

Any reader familiar with linear dynamical systems will notice a strong similarity between this model and a discrete-time autonomous LDS. The idea underlying model (3.1) - (3.6) is that the evolution of the frames in a video depicting the face of a person, can be ascribed to the evolution of an underlying latent variable comprising pose, illumination and expression variations. In addition to this, a constant latent variable related to identity influences the aspect of each frame, while remaining the same across all videos of one

---

Table 3.1: Notation

---

$\mathbf{x}_{ij}^t$	–	Observation relative to the $t$ -th frame from the $j$ -th video of the $i$ -th individual.
$\epsilon_{ij}^t$	–	Observation noise affecting the same frame.
$\mathbf{w}_{ij}^t$	–	Private latent variable underlying the observation.
$\mathbf{v}_{ij}^t$	–	Process noise affecting the evolution of the private latent variable through time.
$\mathbf{h}_i$	–	Public latent variable shared across all videos of person $i$ .
$\mathbf{z}_{ij}^t$	–	Concatenation of private and public latent variable underlying a specific frame: $\mathbf{z}_{ij}^t = [\mathbf{h}_i \mathbf{w}_{ij}^t]^T$ .
$\boldsymbol{\mu}$	–	Mean value, averaged across all frames of every video of each person.
$\mathbf{x}_{ij}$	–	Column vector containing $j$ -th video of person $i$ , concatenating all frames in a single vector.
$\mathbf{w}_{ij}$	–	Column vector comprising the private latent variables associated to the $j$ -th video of person $i$ .
$\mathbf{x}_i^t$	–	Column vector containing the $t$ -th frame from all the videos of person $i$ , concatenated in a single vector.
$\mathbf{w}_i^t$	–	Column vector grouping all the private latent variables of person $i$ at time $t$ .
$\mathbf{z}_i^t$	–	Column vector grouping all the latent variables, both private and public, for person $i$ at time $t$ .
$\mathbf{x}_i$	–	Column vector with all the frames from all the videos available for person $i$ .
$\mathbf{z}_i$	–	Column vector grouping all the latent variables, both private and public, for person $i$ .
$\mathbf{X}$	–	Column vector assembling the totality of the observations.
$\mathbf{w}$	–	Column vector with all the underlying private latent variables.
$\mathbf{h}$	–	Column vector with all the underlying public latent variables.
$\mathbf{Z}$	–	Column vector grouping the totality of the latent variables.
$f$	–	Number of dimensions in the observation space.
$D_w$	–	Number of dimensions for the private latent space.
$D_h$	–	Number of dimensions for the public latent space.
$I$	–	Number of different individuals.
$J$	–	Number of videos per individual.
$T$	–	Number of frames per video.

---

person. Again, note that this model is an artificial description and we don't claim it has any adherence with the physics behind the shooting of a video. However, just like PLDA, it exhibits a good discriminatory power and an intuitive explanation for each parameter while retaining mathematical tractability in terms of learning and inference, to which we will now turn our attention.

### 3.2 Learning and recognition

The learning phase, in analogy with PLDA, requires the estimation of parameter  $\theta = (\mu, \mathbf{F}, \mathbf{G}, \Sigma, \mathbf{A}_{11, \dots, IJ})$ . Here matrices  $\mathbf{A}_{ij}$  represent a novelty specific to dynamic PLDA and they are responsible for the evolution of the private latent variables. For  $\mathbf{w}_{ij}^t$  we expect slow variations not prone to divergence. Hence we anticipate the dominant eigenvalues of  $\mathbf{A}_{ij}$  to be positive and slightly less than 1.

As mentioned in section 2.4, the first step to perform learning and later recognition is to consider the joint likelihood of observations and latent variables. To avoid complicating an already convoluted notation, in this section we assume that every individual has exactly  $J$  videos in the training set, each consisting of  $T$  frames. Given this premise we have:

$$\begin{aligned}
 p(\mathbf{X}, \mathbf{Z}|\theta) &= \underbrace{\prod_{i=1}^I p(\mathbf{x}_i|\mathbf{z}_i, \theta)}_{P_1} \underbrace{\prod_{i=1}^I p(\mathbf{z}_i|\theta)}_{P_2} \\
 P_1 &= \prod_{i=1}^I \prod_{j=1}^J \prod_{t=1}^T \underbrace{p(\mathbf{x}_{ij}^t|\mathbf{w}_{ij}^t, \mathbf{h}_i, \theta)}_{\mathcal{N}(\mathbf{x}_{ij}^t|\mu + \mathbf{F}\mathbf{h}_i + \mathbf{G}\mathbf{w}_{ij}^t, \Sigma)} \\
 P_2 &= \left( \prod_{i=1}^I \underbrace{p(\mathbf{h}_i|\theta)}_{\mathcal{N}(\mathbf{h}_i|\mathbf{0}, \mathbf{I})} \right) \prod_{i=1}^I \prod_{j=1}^J \underbrace{p(\mathbf{w}_{ij}^1|\theta)}_{\mathcal{N}(\mathbf{w}_{ij}^1|\mathbf{0}, \mathbf{I})} \prod_{t=2}^T \underbrace{p(\mathbf{w}_{ij}^t|\mathbf{w}_{ij}^{t-1}, \theta)}_{\mathcal{N}(\mathbf{w}_{ij}^t|\mathbf{A}_{ij}\mathbf{w}_{ij}^{t-1}, \mathbf{I})}
 \end{aligned}$$

$$\begin{aligned}
 \ln p(\mathbf{X}, \mathbf{Z}|\theta) &= \ln P_1 + \ln P_2 \\
 &= \sum_{ijt} \ln \mathcal{N}(\mathbf{x}_{ij}^t|\mu + \mathbf{F}\mathbf{h}_i + \mathbf{G}\mathbf{w}_{ij}^t, \Sigma) + \\
 &\quad \sum_i \ln \mathcal{N}(\mathbf{h}_i|\mathbf{0}, \mathbf{I}) + \sum_{ij} \ln \mathcal{N}(\mathbf{w}_{ij}^1|\mathbf{0}, \mathbf{I}) + \\
 &\quad \sum_{ij} \sum_{t=2}^T \ln \mathcal{N}(\mathbf{w}_{ij}^t|\mathbf{A}_{ij}\mathbf{w}_{ij}^{t-1}, \mathbf{I}) \\
 &= \sum_{ijt} \left\{ \ln \left( \frac{1}{(2\pi)^{f/2} |\Sigma|^{1/2}} \right) + \right. \\
 &\quad \left. - \frac{1}{2} (\mathbf{x}_{ij}^t - \mu + \mathbf{F}\mathbf{h}_i + \mathbf{G}\mathbf{w}_{ij}^t)^T \Sigma^{-1} (\mathbf{x}_{ij}^t - \mu + \mathbf{F}\mathbf{h}_i + \mathbf{G}\mathbf{w}_{ij}^t) \right\} + \\
 &\quad \sum_i \left\{ \ln \left( \frac{1}{(2\pi)^{D_h/2}} \right) - \frac{1}{2} \mathbf{h}_i^T \mathbf{h}_i \right\} + \sum_{ij} \left\{ \ln \left( \frac{1}{(2\pi)^{D_w/2}} \right) - \frac{1}{2} \mathbf{w}_{ij}^{1T} \mathbf{w}_{ij}^1 \right\} + \\
 &\quad \sum_{ij} \sum_{t=2}^T \left\{ \ln \left( \frac{1}{(2\pi)^{D_w/2}} \right) - \frac{1}{2} (\mathbf{w}_{ij}^t - \mathbf{A}_{ij}\mathbf{w}_{ij}^{t-1})^T (\mathbf{w}_{ij}^t - \mathbf{A}_{ij}\mathbf{w}_{ij}^{t-1}) \right\} \quad (3.7)
 \end{aligned}$$

Just like for PLDA, this procedure is a prerequisite for carrying out learning through the EM algorithm. It is based on splitting the joint likelihood into a prior on the latent variables and a conditional for the observations, merely exploiting Bayes' law. Each of them is further decomposed into the elementary building blocks of model (3.1) - (3.6). These are then additively separated thanks to the logarithmic function, which also allows to further simplify the expression by inverting the exponential from the gaussian distributions. Note that, unsurprisingly, PLDA can be obtained as a special case of DPLDA by choosing  $T = 1$ , that is to say by considering single frames as videos of unitary length.

### 3.2.1 M step

This phase is very similar to its analogous for PLDA. The main difference is the need to estimate matrix  $A$  and to take into account temporal dependencies. In order to evaluate the parameters  $\theta = (\mu, F, G, \Sigma, A_{ij})$ , we take the expected value of (3.7) with respect to the PDF  $p(\mathbf{h}_i \mathbf{w}_{ij}^t | x)$ :

$$\begin{aligned} E[\ln p(\mathbf{X}, \mathbf{Z} | \theta)] = & - \sum_{ijt} \left\{ \frac{1}{2} \ln |\Sigma| + \frac{1}{2} (\mathbf{x}_{ij}^t - \mu)^T \Sigma^{-1} (\mathbf{x}_{ij}^t - \mu) + \frac{1}{2} E[\mathbf{z}_{ij}^{tT} \mathbf{B}^T \Sigma^{-1} \mathbf{B} \mathbf{z}_{ij}^t] + \right. \\ & \left. - E[\mathbf{z}_{ij}^{tT} \mathbf{B}^T \Sigma^{-1} (\mathbf{x}_{ij}^t - \mu)] \right\} \\ & - \sum_{ij} \sum_{t=2}^T \left\{ \frac{1}{2} E[\mathbf{w}_{ij}^{t-1T} \mathbf{A}_{ij}^T \mathbf{A}_{ij} \mathbf{w}_{ij}^{t-1}] - E[\mathbf{w}_{ij}^{t-1T} \mathbf{A}_{ij}^T \mathbf{w}_{ij}^t] \right\} + K, \end{aligned} \quad (3.8)$$

where  $K$  includes all the terms independent of  $\theta$  and we have defined

$$\begin{aligned} \mathbf{B} &= [\mathbf{F} \quad \mathbf{G}], \\ \mathbf{z}_{ij}^t &= \begin{bmatrix} \mathbf{h}_i \\ \mathbf{w}_{ij}^t \end{bmatrix}. \end{aligned}$$

Just like in section 2 we compute the derivatives of expression (3.8) with respect to each parameter. By setting them to zero we have:

$$\mu = \frac{1}{IJT} \sum_{ijt} \mathbf{x}_{ij}^t, \quad (3.9)$$

$$\mathbf{B} = \left( \sum_{ijt} (\mathbf{x}_{ij}^t - \mu) E[\mathbf{z}_{ij}^{tT}] \right) \left( \sum_{ijt} E[\mathbf{z}_{ij}^t \mathbf{z}_{ij}^{tT}] \right)^{-1}, \quad (3.10)$$

$$\Sigma = \frac{1}{IJT} \sum_{ijt} \text{diag} \{ (\mathbf{x}_{ij}^t - \mu)(\mathbf{x}_{ij}^t - \mu)^T - 2\mathbf{B} E[\mathbf{z}_{ij}^t] (\mathbf{x}_{ij}^t - \mu)^T + \mathbf{B} E[\mathbf{z}_{ij}^t \mathbf{z}_{ij}^{tT}] \mathbf{B}^T \}. \quad (3.11)$$

Note the complete analogy with equations (2.33) - (2.35), the only difference laying in the summation over index  $t$ . Maximization of (3.8) with respect to  $A_{ij}$  yields

$$\begin{aligned} \frac{\partial}{\partial \mathbf{A}_{ij}} \left\{ - \sum_{ij} \sum_{t=2}^T \frac{1}{2} \text{tr} [\mathbf{A}_{ij}^T \mathbf{A}_{ij} \text{Var}(\mathbf{w}_{ij}^{t-1})] + \frac{1}{2} E[\mathbf{w}_{ij}^{t-1}]^T \mathbf{A}_{ij}^T \mathbf{A}_{ij} E[\mathbf{w}_{ij}^{t-1}] \right. \\ \left. - \text{tr} [\mathbf{A}_{ij} \text{Cov}(\mathbf{w}_{ij}^{t-1}, \mathbf{w}_{ij}^t)] - E[\mathbf{w}_{ij}^{t-1}]^T \mathbf{A}_{ij}^T E[\mathbf{w}_{ij}^t] \right\} = 0, \end{aligned}$$

where we have used

$$E[\mathbf{x}^T \mathbf{A} \mathbf{y}] = \text{tr}[\mathbf{A}^T \text{Cov}(\mathbf{x}, \mathbf{y})] + E(\mathbf{x})^T \mathbf{A} E(\mathbf{y}). \quad (3.12)$$

Hence

$$\begin{aligned} \sum_{t=2}^T \left\{ \mathbf{A}_{ij} \text{Var}(\mathbf{w}_{ij}^{t-1}) - \text{Cov}(\mathbf{w}_{ij}^{t-1}, \mathbf{w}_{ij}^t)^T + E[\mathbf{w}_{ij}^t] E[\mathbf{w}_{ij}^{t-1}]^T (\mathbf{A}_{ij} - \mathbf{I}) \right\} &= 0 \\ \mathbf{A}_{ij} &= \left( \sum_{t=2}^T E[\mathbf{w}_{ij}^t \mathbf{w}_{ij}^{t-1 T}] \right) \left( \sum_{t=2}^T E[\mathbf{w}_{ij}^{t-1} \mathbf{w}_{ij}^{t-1 T}] \right)^{-1} \end{aligned} \quad (3.13)$$

As a special case, if we imposed the same state evolution matrix for each  $\mathbf{w}_{ij}$ , that is  $\mathbf{A}_{ij} = \mathbf{A}$ , we would get

$$\mathbf{A} = \left( \sum_{ij} \sum_{t=2}^T E[\mathbf{w}_{ij}^t \mathbf{w}_{ij}^{t-1 T}] \right) \left( \sum_{ij} \sum_{t=2}^T E[\mathbf{w}_{ij}^{t-1} \mathbf{w}_{ij}^{t-1 T}] \right)^{-1}. \quad (3.14)$$

### 3.2.2 E step

In this phase we need to compute the first and second order moments in formulas (3.9), (3.10), (3.11) and (3.13), where the expectation is conditional on all the available observations  $\mathbf{X}$ . As already pointed out, equations (3.1) - (3.6) describe a system largely similar to an autonomous LDS. As  $i$  and  $j$  vary, a number  $IJ$  of such systems is spanned. The main difference with a classical LDS, however, is the presence of a latent variable  $h_i$  shared across all the  $J$  videos of person  $i$ . The most intuitive solution to this problem is to combine the observations and the latent variables for each video in analogy with (2.25), resulting in the augmented system

$$\begin{bmatrix} \mathbf{x}_{i1}^t \\ \mathbf{x}_{i2}^t \\ \vdots \\ \mathbf{x}_{iJ}^t \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu} \\ \vdots \\ \boldsymbol{\mu} \end{bmatrix} + \begin{bmatrix} \mathbf{F} & \mathbf{G} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{G} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{F} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{G} \end{bmatrix} \begin{bmatrix} \mathbf{h}_i \\ \mathbf{w}_{i1}^t \\ \vdots \\ \mathbf{w}_{iJ}^t \end{bmatrix} + \begin{bmatrix} \boldsymbol{\epsilon}_{i1}^t \\ \boldsymbol{\epsilon}_{i2}^t \\ \vdots \\ \boldsymbol{\epsilon}_{iJ}^t \end{bmatrix}, \quad (3.15)$$

$$\begin{bmatrix} \mathbf{h}_i \\ \mathbf{w}_{i1}^{t-1} \\ \vdots \\ \mathbf{w}_{iJ}^{t-1} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & & & \\ & \mathbf{A}_{i1} & & \\ & & \mathbf{A}_{i2} & \\ & & & \ddots \\ & & & & \mathbf{A}_{iJ} \end{bmatrix} \begin{bmatrix} \mathbf{h}_i \\ \mathbf{w}_{i1}^{t-1} \\ \vdots \\ \mathbf{w}_{iJ}^{t-1} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \mathbf{v}_{i1}^t \\ \vdots \\ \mathbf{v}_{iJ}^t \end{bmatrix}, \quad (3.16)$$

which can be re-written in a more compact form as

$$\begin{cases} \mathbf{x}_i^t = \bar{\mathbf{C}} \mathbf{z}_i^t + \bar{\boldsymbol{\mu}} + \boldsymbol{\epsilon}_i^t \\ \mathbf{z}_i^t = \bar{\mathbf{A}}_i \mathbf{z}_i^{t-1} + \mathbf{v}_i^t \end{cases}, \quad (3.17)$$

where

$$\boldsymbol{\epsilon}_i^t \sim \mathcal{N}(\mathbf{0}, \bar{\boldsymbol{\Sigma}}), \quad \bar{\boldsymbol{\Sigma}} = \begin{bmatrix} \boldsymbol{\Sigma} & & & \\ & \boldsymbol{\Sigma} & & \\ & & \ddots & \\ & & & \boldsymbol{\Sigma} \end{bmatrix} \quad (3.18)$$

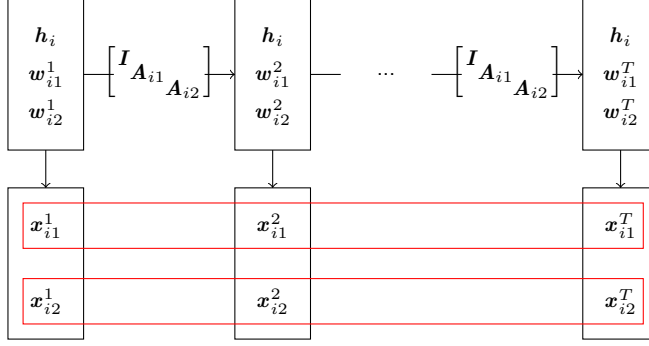


Figure 3.1: Schematic for the LDS. The arrows depict the flow of information as each latent variable is realized and influences both its future value and the current state. The frames from the same video are surrounded by red rectangles.

$$\mathbf{v}_i^t \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}), \quad \mathbf{\Gamma} = \begin{bmatrix} \mathbf{0} & & & \\ & \mathbf{I} & & \\ & & \ddots & \\ & & & \mathbf{I} \end{bmatrix} \quad (3.19)$$

and

$$\begin{aligned} \mathbf{x}_i^t, \bar{\boldsymbol{\mu}}, \boldsymbol{\epsilon}_i^t &\in \mathbb{R}^{Jfx1}, \\ \mathbf{z}_i^t, \mathbf{v}_i^t &\in \mathbb{R}^{D_h + JD_w}, \\ \bar{\mathbf{C}} &\in \mathbb{R}^{Jfx(D_h + JD_w)}, \\ \bar{\mathbf{A}}_i &\in \mathbb{R}^{(D_h + JD_w) \times (D_h + JD_w)}. \end{aligned}$$

This way we successfully put system (3.17) in a canonical LDS form, while guaranteeing that the original properties desired for the model are retained. In this respect, note that that we are forcing the identity variable  $\mathbf{h}_i$  to be shared for all the videos of person  $i$  while remaining the same across all frames. The private variable  $\mathbf{w}_{ij}^t$  is instead subject to the typical dynamics of a linear system. Figure 3.1 depicts a schematic for the model with  $J = 2$ . Here, much like in a graphical model, the arrows provide information about the conditional dependencies in effect. The red rectangles enclose all the frames from a single video.

It should be observed that reducing our model to a canonical LDS came at the cost of multiplying the size of the observations and the hidden state by a factor of  $J$ , which results in an undesirable computational complexity for Kalman filtering, increasing with the cube of  $J$ . In section 3.3.1 we provide a less intuitive yet equivalent formulation which circumvents this issue. Hereafter, borrowing from the control systems literature, we will use the word “state” as a synonym for latent variable.

The advantage of dealing with a canonical LDS is the chance of resorting to the related, well established estimate techniques, epitomized by the Kalman filter. In the presence of jointly gaussian observations and latent variables, the Kalman filter gives unbiased conditional estimates for the state while keeping at a minimum the mean square of the prediction error. Also, it scales well with  $T$  since its computational complexity is linear with respect to the length of the filtered sequence. For a thorough introduction to Kalman filtering, [24] and [6] are helpful references. From the latter we borrow the notational



convention, in recalling the filtering equations:

$$\mathbf{m}_i^t = \bar{\mathbf{A}}_i \mathbf{m}_i^{t-1} + \mathbf{K}_i^t (\mathbf{x}_i^t - \bar{\mathbf{C}} \bar{\mathbf{A}}_i \mathbf{m}_i^{t-1}), \quad (3.20)$$

$$\mathbf{V}_i^t = (\mathbf{I} - \mathbf{K}_i^t \bar{\mathbf{C}}) \mathbf{P}_i^{t-1}, \quad (3.21)$$

$$\mathbf{P}_i^{t-1} = \bar{\mathbf{A}}_i \mathbf{V}_i^{t-1} \bar{\mathbf{A}}^T + \mathbf{\Gamma}, \quad (3.22)$$

$$\mathbf{K}_i^t = \mathbf{P}_i^{t-1} \bar{\mathbf{C}}^T (\mathbf{C} \mathbf{P}_i^{t-1} \mathbf{C}^T + \mathbf{\Sigma})^{-1}, \quad (3.23)$$

with first update

$$\mathbf{m}_i^1 = \mathbf{m}_i^0 + \mathbf{K}_i^1 (\mathbf{x}_i^1 - \bar{\mathbf{C}} \mathbf{m}_i^0), \quad (3.24)$$

$$\mathbf{V}_i^1 = (\mathbf{I} - \mathbf{K}_i^1 \bar{\mathbf{C}}) \mathbf{P}_i^0, \quad (3.25)$$

and initial conditions, resulting from priors (3.5) - (3.6), given by

$$\mathbf{m}_i^0 = [\mathbf{0}, \mathbf{0}, \dots, \mathbf{0}]^T, \quad (3.26)$$

$$\mathbf{P}_i^0 = \begin{bmatrix} \mathbf{I} & & & \\ & \mathbf{I} & & \\ & & \ddots & \\ & & & \mathbf{I} \end{bmatrix}. \quad (3.27)$$

These equations allow us to compute the filtering probabilities

$$P(\mathbf{z}_i^t | \mathbf{x}_i^1, \dots, \mathbf{x}_i^t) = \mathcal{N}(\mathbf{m}_i^t, \mathbf{V}_i^t), \quad (3.28)$$

and the conditionals

$$c_i^t = P(\mathbf{x}_i^t | \mathbf{x}_i^1, \dots, \mathbf{x}_i^{t-1}) = \mathcal{N}(\mathbf{x}_i^t | \bar{\mathbf{C}} \bar{\mathbf{A}}_i \mathbf{m}_i^{t-1}, \mathbf{\Sigma} + \bar{\mathbf{C}} \mathbf{P}_i^{t-1} \bar{\mathbf{C}}^T), \quad (3.29)$$

$$c_i^1 = P(\mathbf{x}_i^1) = \mathcal{N}(\mathbf{x}_i^1 | \bar{\mathbf{C}} \mathbf{m}_i^0, \mathbf{\Sigma} + \bar{\mathbf{C}} \mathbf{P}_i^0 \bar{\mathbf{C}}^T) \quad (3.30)$$

Note that the probability (3.28) is for  $\mathbf{z}_i^t$  conditional on the observations up to the same time step  $t$ , whereas for the E step we want to condition on all the available information  $\mathbf{X}$  up to the final time step  $T$ . To take into account the totality of the information, the filtering phase is followed by Kalman smoothing, which is a backward propagation of information that can be formulated as follows:

$$\widehat{\mathbf{m}}_i^T = \mathbf{m}_i^T, \quad (3.31)$$

$$\widehat{\mathbf{V}}_i^T = \mathbf{V}_i^T, \quad (3.32)$$

$$\widehat{\mathbf{m}}_i^t = \mathbf{m}_i^t + \mathbf{J}_i^t (\widehat{\mathbf{m}}_i^{t+1} - \bar{\mathbf{A}}_i \mathbf{m}_i^t), \quad (3.33)$$

$$\widehat{\mathbf{V}}_i^t = \mathbf{V}_i^t + \mathbf{J}_i^t (\widehat{\mathbf{V}}_i^{t+1} - \mathbf{P}_i^t) \mathbf{J}_i^{tT}, \quad (3.34)$$

$$\mathbf{J}_i^t = \mathbf{V}_i^t \bar{\mathbf{A}}_i^T (\mathbf{P}_i^t)^{-1}, \quad (3.35)$$

yielding the smoothing probabilities

$$P(\mathbf{z}_i^t | \mathbf{X}) = \mathcal{N}(\mathbf{z}_i^t | \widehat{\mathbf{m}}_i^t, \widehat{\mathbf{V}}_i^t), \quad (3.36)$$

and second order moments

$$\text{Cov} [\mathbf{z}_i^{t-1}, \mathbf{z}_i^t | \mathbf{X}] = \mathbf{J}_i^{t-1} \widehat{\mathbf{V}}_i^t. \quad (3.37)$$

Remembering the definition of  $\mathbf{z}_i^t = [\mathbf{h}_i, \mathbf{w}_{i1}^t, \dots, \mathbf{w}_{iJ}^t]^T$ , from the first and second order moments in formulas (3.33), (3.34), (3.36), we can extract the sufficient statistics necessary for the M step.

### 3.2.3 Recognition

As in chapter 2, recognition problems can be solved, provided we have an algorithm to estimate the likelihood that  $J$  videos share the same identity. Namely we have to estimate  $P(\mathbf{x}_i) = P(\mathbf{x}_i^1, \dots, \mathbf{x}_i^T)$ , where  $\mathbf{x}_i^t$  is defined according to (3.17) as the column vector containing all the pixels of all videos of person  $i$  at time  $t$ . Resorting to equation (3.29), said likelihood can be computed through Kalman filtering as

$$P(\mathbf{x}_i) = \prod_{t=1}^T c_i^t. \quad (3.38)$$

## 3.3 Implementation Details

### 3.3.1 Video concatenation

As mentioned in the previous section, augmenting the state and the observation vectors according to model (3.15) becomes more and more inefficient as the number of gallery videos per individual  $J$  increases. Typically, the bottleneck for Kalman filters is represented by the inversion of matrix  $(\mathbf{C}\mathbf{P}_i^{t-1}\mathbf{C}^T + \mathbf{\Sigma})$  in equation (3.23), which for the above mentioned model would be of size  $Jf \times Jf$ . Its inversion would require a number of operations proportional to the cube of  $Jf$ . A more desirable model would allow us to transfer across videos the information learnt about the identity latent variable  $\mathbf{h}_i$ , while maintaining a linear complexity with respect to  $J$ . Another limitation of the above formulation is that it does not easily generalize to the case of videos of different length.

Consider the model in figure 3.2, alternative to the one previously discussed in figure 3.1. According to this new setting, depicting the case  $J = 2$ , videos of the same person can be joined end-to-front, avoiding the inefficient state and observation augmentation. This is shown by the rectangles in red, which again delimit the frames within the same video. Exploiting this idea, frames from videos of the same person are concatenated in one single macro-video of length  $JT$ . Crucially, in correspondence to the transition from one video to the next, the matrix responsible for the evolution of the latent variables changes. More specifically, its lower diagonal block, the one responsible for the evolution of the private latent variable, is equal to  $\mathbf{A}_{i1}$  for the whole duration of the first video. When the video ends, it becomes zero for one iteration, making sure that the prior for the private latent variable  $\mathbf{w}_{i2}^1$  is set to a gaussian with zero mean and unit covariance, blocking the transfer of information from video to video. For the next iteration, i.e. for the transition from the first to the second frame of video 2, the lower diagonal block takes the value  $\mathbf{A}_{i2}$ , and retains it for the whole duration of the second video. If a third video was available, another

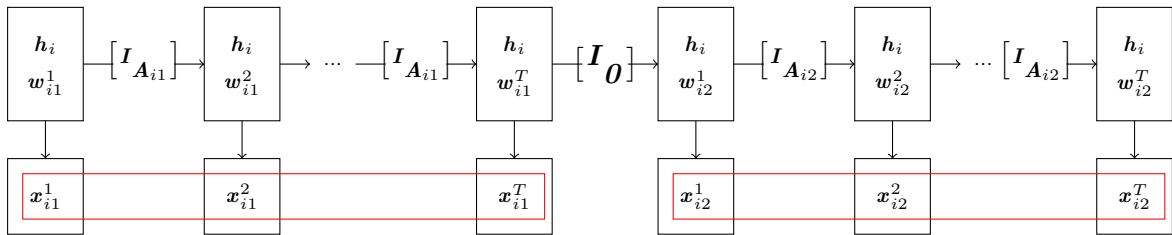


Figure 3.2: Schematic for the LDS. Videos of the same person are joined end-to-front in a single macro video. In correspondence of a transition from a video to the next one, the state evolution matrix changes.

“zero transition” would be present between video 2 and video 3. At the same time, the upper diagonal block, responsible for the evolution of the public latent variable  $\mathbf{h}_i$ , remains fixed at  $\mathbf{I}$ . This effectively forces  $\mathbf{h}_i$  to be constant across all videos and all frames of the same individual and allows us to transfer the related information across videos.

To show the equivalence between this formulation and the one featuring the augmented state, we have to prove that the time-varying model resulting from the different state transition matrices described above, applied to the macro-video, enforces the desired priors  $\mathbf{w}_{ij}^1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . For  $j = 1$  this is trivial, since  $\mathbf{w}_{ij}^1$  underlays the first frame of the macro-video, for which such prior is imposed as usual. For  $j > 1$ , model (3.17)-(3.19) still holds, after redefining

$$\bar{\mathbf{C}} = [\mathbf{F}, \mathbf{G}], \quad (3.39)$$

$$\bar{\mathbf{\Gamma}} = \begin{bmatrix} \mathbf{0} & \\ & \mathbf{I} \end{bmatrix} \quad (3.40)$$

$$\mathbf{z}_i^t = [\mathbf{h}_i, \mathbf{w}_i^t]^T, \quad (3.41)$$

$$\mathbf{w}_i^t = \mathbf{w}_{ij}^k, \quad (3.42)$$

$$\mathbf{x}_i^t = \mathbf{x}_{ij}^k, \quad (3.43)$$

where  $j, k$  are defined as

$$j = (t - 1)/T + 1,$$

$$k = (t - 1) \% T + 1.$$

Here ‘/’ denotes the integer division and ‘%’ its remainder. Finally, state matrix  $\mathbf{A}_i$  varies as discussed above, according to figure 3.2.

Given said definitions, let us take as an example  $j = 2$ , without any loss of generality. We want to estimate

$$P(\mathbf{z}_i^{T+1} | \mathbf{x}_i^1, \dots, \mathbf{x}_i^{T+1}) = P\left(\begin{bmatrix} \mathbf{h}_i \\ \mathbf{w}_{i2}^1 \end{bmatrix} | \mathbf{x}_{i1}^1, \dots, \mathbf{x}_{i1}^T, \mathbf{x}_{i2}^1\right) = \mathcal{N}(\mathbf{m}_i^{T+1}, \mathbf{V}_i^{T+1}) \quad (3.44)$$

and look into the PDF of its lower half  $\mathbf{w}_{i2}^1$ . Notice that the current state matrix

$$\mathbf{A}_i^{T+1} = \begin{bmatrix} \mathbf{I} & \\ & \mathbf{0} \end{bmatrix} \quad (3.45)$$

directly affects the update of  $\mathbf{P}_i^T$  and  $\mathbf{m}_i^{T+1}$  as follows

$$\mathbf{P}_i^T = \begin{bmatrix} \mathbf{I} & \\ & \mathbf{0} \end{bmatrix} \mathbf{V}_i^T \begin{bmatrix} \mathbf{I} & \\ & \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \\ & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \mathbf{V}_{i,h}^T & \\ & \mathbf{I} \end{bmatrix} \quad (3.46)$$

$$\begin{aligned} \mathbf{m}_i^{T+1} &= \begin{bmatrix} \mathbf{I} & \\ & \mathbf{0} \end{bmatrix} \mathbf{m}_i^T + \mathbf{K}_i^{T+1} \left( \mathbf{x}_i^{T+1} - [\mathbf{F} \mathbf{G}] \begin{bmatrix} \mathbf{I} & \\ & \mathbf{0} \end{bmatrix} \mathbf{m}_i^T \right) \\ &= \begin{bmatrix} \bar{\mathbf{h}}_i^T \\ \mathbf{0} \end{bmatrix} + \mathbf{K}_i^{T+1} \left( \mathbf{x}_i^{T+1} - \begin{bmatrix} \mathbf{F} \bar{\mathbf{h}}_i^T \\ \mathbf{0} \end{bmatrix} \right), \end{aligned} \quad (3.47)$$

where we have called  $\bar{\mathbf{h}}_i^T$ ,  $\mathbf{V}_{i,h}^T$  the estimated mean and covariance for  $\mathbf{h}_i$  conditional on the information up to time step  $T$ . Comparing these results with (3.24) - (3.27), it should be clear that this model behaves exactly as desired. That is to say, in correspondence to

a transition from the first video to the second, it behaves like a freshly started Kalman filter with prior beliefs,

$$\mathbf{P}_i^0 = \begin{bmatrix} \mathbf{V}_{i,h}^T & \\ & \mathbf{I} \end{bmatrix}, \quad (3.48)$$

$$\mathbf{m}_i^0 = \begin{bmatrix} \bar{\mathbf{h}}_i^T \\ \mathbf{0} \end{bmatrix}, \quad (3.49)$$

thus allowing for an unaltered transfer of information about  $\mathbf{h}_i$  while resetting the information about  $\mathbf{w}_{i2}$  just like a prior  $\mathbf{w}_{i2}^1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Equivalently it can be said that, under the model just discussed,

$$P(\mathbf{w}_{i2}^1 | \mathbf{x}_{i1}^1, \dots, \mathbf{x}_{i1}^T, \mathbf{x}_{i2}^1) = P(\mathbf{w}_{i2}^1 | \mathbf{x}_{i2}^1). \quad (3.50)$$

### 3.3.2 Offline computations

By inspecting the filtering and smoothing equations (3.20) - (3.23) and (3.33) - (3.35), one should notice that the heaviest computations can be carried out offline, since the measurements  $\mathbf{X}$  only affect the estimates of the expectations. This is especially advantageous in the case of a unique state matrix shared among all videos, that is  $\mathbf{A}_{ij} = \mathbf{A}$ . In this case, the most demanding operation, the inversion of matrices  $(\mathbf{C}\mathbf{P}_i^{t-1}\mathbf{C}^T + \mathbf{\Sigma})$ , can be carried out once and for all videos preliminarily, thus significantly alleviating the computational burden associated with DPLDA.

In addition, the inversion of matrix  $(\mathbf{C}\mathbf{P}_i^{t-1}\mathbf{C}^T + \mathbf{\Sigma})$  can be carried out efficiently exploiting Woodbury's formula

$$(\mathbf{C}\mathbf{P}_i^{t-1}\mathbf{C}^T + \mathbf{\Sigma})^{-1} = \mathbf{\Sigma}^{-1} - \mathbf{\Sigma}^{-1}(\mathbf{C}(\mathbf{C}^T\mathbf{\Sigma}^{-1}\mathbf{C} + [\mathbf{P}_i^{t-1}]^{-1})^{-1}\mathbf{C}^T)\mathbf{\Sigma}^{-1}, \quad (3.51)$$

wherein the inversion of  $\mathbf{\Sigma}$  is simple to compute due to its diagonal structure.

### 3.3.3 Inference

In order to decide the most likely identity for a probe video, in principle we can estimate the likelihood that it shares its identity with all the subjects in the gallery, according to (3.38) and pick the person that maximizes it. This would imply concatenating the probe video to all the gallery videos of a subject as discussed in subsection 3.3.1, and performing Kalman filtering on this macro video. Suppose that the probe video consists of  $T$  frames and that the gallery is made up of  $I_{gal}$  individuals each of which is depicted in  $J$  different videos of length  $T$ . To try out all the possible combinations one would have to run  $I_{gal}$  Kalman filters of length  $(J+1)T$ .

A more efficient approach consists in running offline a Kalman filter of length  $JT$  on each gallery subject, consequently storing the learnt values for the public latent variable. As soon as a probe video is available for identification,  $I_{gal}$  Kalman filters of length  $T$  can be run along its frames utilizing a different initial state estimate  $\mathbf{m}_i^0$  in each of the runs, corresponding to the values learnt for  $\mathbf{h}_i$  from each different gallery individual which had previously been stored. Supposing  $I_{probe}$  different videos need to be identified, this simple trick reduces the computational requirements from  $O(I_{gal}I_{probe}(J+1)T)$  to  $O(I_{gal}JT + I_{probe}I_{gal}T)$ .

From a probabilistic perspective this approach basically consists in recognizing that the joint probability of gallery and probe videos can be split as

$$P(\mathbf{x}_{i1} \dots \mathbf{x}_{iJ} \mathbf{x}_{probe}) = P(\mathbf{x}_{i1} \dots \mathbf{x}_{iJ})P(\mathbf{x}_{probe} | \mathbf{x}_{i1} \dots \mathbf{x}_{iJ}) \quad (3.52)$$

into a marginal probability for the gallery videos alone and a conditional probability for the probe based on the information learnt from the gallery videos. These 2 probabilities can be assessed using two separate Kalman filters.

Note that the evaluation of said probabilities entails the evaluation of a multivariate gaussian PDF, thus requiring the computation of the determinant for matrix  $(\mathbf{C}\mathbf{P}_i^{t-1}\mathbf{C}^T + \mathbf{\Sigma})$ . Consider for instance a video consisting of 90x90 pixel frames, which can be concatenated in a single feature vector of size  $f = 8100$ . The covariance matrices  $(\mathbf{C}\mathbf{P}_i^{t-1}\mathbf{C}^T + \mathbf{\Sigma})$  will thus have size 8100x8100. With double precision, if the geometric mean of its eigenvalues falls outside the interval  $[0.91, 1.09]$ , direct evaluation of the matrix's determinant will result in underflow or overflow. For this reason, said computation is carried out in the domain of logarithms. By exploiting the positive definiteness of the covariance matrix, we can resort to a Cholesky decomposition to obtain its eigenvalues and compute the logarithm of its determinant  $l$  as follows. Given the Cholesky decomposition

$$(\mathbf{C}\mathbf{P}_i^{t-1}\mathbf{C}^T + \mathbf{\Sigma}) = \mathbf{U}\mathbf{U}^T,$$

then

$$l = 2 \sum_{k=1}^f \log(\mathbf{U}_{kk}),$$

where  $\mathbf{U}_{kk}$  denotes the  $k$ -th element on the diagonal of matrix  $\mathbf{U}$ . With this simple idea we can circumvent the problem of underflow and overflow in an efficient manner while exploiting the LDSs by direct comparison of different log-likelihoods.



## 4. Pipeline and simulations

The models presented so far are applicable to face images in ideal viewing conditions, wherein faces have been located, cropped and appear under a frontal pose. Furthermore, although model noise is accounted for with a gaussian term, it is unrealistic to assume that disturbances such as occlusions and different illuminations will always conform to our statistical convenience. For this reason, we introduce a recognition pipeline wherein DPLDA represents the third and last stage. In order to obtain a good input for DPLDA, we employ a first stage for face detection and frontalization, followed by one of feature extraction. The first stage makes sure that pixel locations can be consistently treated as specific to the same area of the face, e.g. guaranteeing that pixel (40,32) always lays on the tip of the nose. The second stage extracts from each frontalized image a representation that is more suitable for face recognition, guaranteeing a higher robustness margin to pixel noise. Finally, DPLDA takes care of the recognition task, with improved performance thanks to the high quality features resulting from the combination of the two previous stages.

### 4.1 Face frontalization

The problems of face detection and alignment are an essential prerequisite for face recognition, guaranteeing the subsequent steps of a recognition pipeline to work on pixel values consistently related to the same areas of the face. For this purpose we adopt the FAR framework, as detailed in [25]. This technique performs joint face landmark detection and frontalization reconstruction, employing a limited amount of training data from frontal-pose faces, without any need to resort to popular yet computationally expensive 3D face models.

A pivotal idea for this technique is the consideration that the rank for a frontal image (expressed as a matrix) is smaller than the rank from the same face in any other pose, due to the natural symmetry of the human face. Specifically, the problem of face frontalization can be viewed as that of finding the warping of a random-pose face image that gives the closest result to a frontal image of the same face. In this respect, we define  $\mathbf{x}(\mathbf{p}) = \mathbf{x}(\mathcal{W}(\mathbf{p}))$  the warping of an image, dependent on warping coordinates  $\mathbf{p}$  and denote by  $\mathbf{X} = \mathcal{R}_{m \times n}(\mathbf{x})$  the reshape operator that acts on a vector yielding an  $n$  by  $m$  matrix.

The idea of a low rank representation for a frontal pose face is then coupled with a belief that only a limited amount of pixels will be corrupted by non-gaussian noise, giving rise to a sparsity constraint on the error matrix. Overall, a warped corrupted image can be expressed as

$$\mathbf{X}(\mathbf{p}) = \mathbf{L} + \mathbf{E} = \sum_{i=1}^k \mathcal{R}_{m \times n}(\mathbf{u}_i)c_i + \mathbf{E}. \quad (4.1)$$

Here  $\mathbf{E}$  represents a sparse error matrix and  $\mathbf{L} = \sum_{i=1}^k \mathcal{R}_{m \times n}(\mathbf{u}_i)c_i$  is the clean frontal-



Figure 4.1: Frontalized faces. Frames from our database of celebrities downloaded from Youtube.

ization of a face image, laying on the low dimensional subspace spanned by the basis  $\langle \mathbf{u}_i \rangle$ ,  $i = 1 \dots k$ . By imposing the low rank and sparsity constraints the following optimization problem arises:

$$\begin{aligned} \arg \min_{\mathbf{L}, \mathbf{E}, c, \Delta p} \quad & \|\mathbf{L}\|_* + \lambda \|\mathbf{E}\|_1 \\ \text{s.t.} \quad & \mathbf{X}(\mathbf{p}) = \mathbf{L} + \mathbf{E}, \quad \mathbf{L} = \sum_{i=1}^k \mathcal{R}_{m \times n}(\mathbf{u}_i) c_i. \end{aligned} \quad (4.2)$$

Here  $\|\cdot\|_*$  denotes the nuclear norm and  $\|\cdot\|_1$  the  $l_1$  norm, which are employed as convex approximations of rank and cardinality respectively. Optimization problem (4.2) is far from trivial, mainly due to its non-linearity with respect to  $\mathbf{p}$  and is solved with an augmented lagrangian. Figure 4.1 shows the result of the frontalization process on a private in-the-wild database we assembled from Youtube.

## 4.2 Feature extraction

So far we dealt with (grey scale) images in their “natural” representation, i.e. as a sequence of pixel values ranging from a minimum (black) to a maximum (white). However, this representation is not robust to affine transformation, differences in illumination, occlusions and, in general, noise. Feature extraction is the branch of machine learning that deals with the research of ways to describe the data, that prove to be particularly suitable for a task of interest, in environments with different extents of non-ideality. In our case, we would like to extract features from images which allow us a good performance for tasks of face recognition in videos, exhibiting robustness to non-gaussian noise. An increasing amount of alternative representations have been proposed and studied in computer vision. Some of them are directly derived from pixel values, such as Local Binary Patterns, others act on different domains, such as Histogram of Oriented Gradients (HOG) or the popular Scale-Invariant Feature Transform (SIFT) both requiring the computation of image gradients.

As we test our algorithm on videos from different and challenging databases, we expect the respective frames to be far from ideal. On the contrary, significant pixel noise is encountered and has to be dealt with by our identification procedure. For this reason we introduce this stage within our pipeline, wherein we choose a robust feature representation for the frontalized images.



Image Gradient Orientation (IGO) methods, as described in [31], represent an interesting alternative to pixel-based features and were shown to outperform other popular descriptors, such as LBP and Gabor features, in tasks of face recognition. In this section we describe the intuition behind it and sketch a proof of its suitability for face recognition in pre-aligned images.

Consider a set of images  $\{\mathbf{I}_i\}$ . Let us compute the respective gradients and indicate by  $\{\boldsymbol{\Phi}_i\}$ , their orientation, which we normalize so that  $\boldsymbol{\Phi}_i \in [0, 2\pi)$ . Taking two different images, we denote their orientation difference as

$$\Delta(\boldsymbol{\Phi}_{ij}) = \boldsymbol{\Phi}_i - \boldsymbol{\Phi}_j. \quad (4.3)$$

Let us define as  $\mathcal{P}$  the indices corresponding to the image support, arranged in lexicographic order and as  $N(\mathcal{P})$  the cardinality of the set. We can then give the following

**Definition** Two images  $\mathbf{I}_i$   $\mathbf{I}_j$  are pixel-wise dissimilar if  $\forall k \in \mathcal{P}$ ,  $\Delta \boldsymbol{\Phi}_{ij}(k) \sim \mathcal{U}[0, 2\pi)$ .

The above condition can be verified with a statistical significance test in order to determine whether 2 images follow the definition or not. In the domain of gradient orientations, let us define the cosine-based correlation between two images as follows:

$$s(\phi_i, \phi_j) = \sum_{k \in \mathcal{P}} \cos[\Delta \boldsymbol{\Phi}_{ij}(k)] = cN(\mathcal{P}), \quad (4.4)$$

where  $c \in [-1, 1]$  is the coefficient providing a measurement of the spatial correlation between images  $\mathbf{I}_i$  and  $\mathbf{I}_j$ . Note for instance that for highly correlated images  $\Delta \boldsymbol{\Phi}_{ij}(k) \approx 0$  hence  $c \rightarrow 1$ .

As we will show, the correlation measurement (4.4) is robust with respect to outliers. In this regard, suppose we can partition the image support  $\mathcal{P}$  into a subset  $\mathcal{P}_2$  of pixels corrupted by outliers (within either image) and a subset  $\mathcal{P}_1$  of ideally noise-less pixels (in both images). A robust correlation measurement between two such images should disregard  $\mathcal{P}_2$  and be well approximated by

$$s_1(\phi_i, \phi_j) = \sum_{k \in \mathcal{P}_1} \cos[\Delta \boldsymbol{\Phi}_{ij}(k)] = c_1 N(\mathcal{P}). \quad (4.5)$$

To show that indeed this property holds, we assume that in  $\mathcal{P}_2$  the images are pixel-wise dissimilar according to the definition given above. This is a reasonable assumption which proves to be well supported by empirical verification and makes the following theorem meaningful for our analysis:

**Theorem** Let  $u(\cdot)$  be a random process and  $u(t) \sim \mathcal{U}[0, 2\pi)$  then:

- $E[\int_{\mathcal{X}} \cos u(t) dt] = 0$  for any non-empty interval  $\mathcal{X} \in \mathbb{R}$ .
- If  $u(\cdot)$  is mean ergodic, then  $\int_{\mathcal{X}} \cos u(t) dt = 0$ .

We also employ the following approximation

$$\int_{\mathcal{X}} \cos[\Delta \boldsymbol{\Phi}_{ij}(t)] dt \simeq \sum_{k \in \mathcal{P}} \cos[\Delta \boldsymbol{\Phi}_{ij}(k)], \quad (4.6)$$

where the integrand in the left hand side is defined in a continuous domain with some notational abuse. By applying the theorem above to the subset of pixels  $\mathcal{P}_2$  corrupted by outliers, we have

$$s_2(\phi_i, \phi_j) = \sum_{k \in \mathcal{P}_2} \cos[\Delta \boldsymbol{\Phi}_{ij}(k)] \simeq 0, \quad (4.7)$$

which in turn yields

$$s(\phi_i, \phi_j) = s_1(\phi_i, \phi_j), \quad (4.8)$$

completing our sketch of proof for the robustness of correlation measurement (4.4) with respect to outliers.

The considerations above allow us to define an interesting distance measure:

$$d^2(\phi_i, \phi_j) = \sum_{k \in \mathcal{P}} \{1 - \cos[\Delta \Phi_{ij}(k)], \} \quad (4.9)$$

which inherits the robustness properties of (4.4). Through basic manipulations, said distance can be rewritten as

$$d^2(\phi_i, \phi_j) = \frac{1}{2} \sum_{k \in \mathcal{P}} \underbrace{\{\cos^2(\phi_i) + \sin^2(\phi_i) + \cos^2(\phi_j) + \sin^2(\phi_j) +}_{2} \quad (4.10)$$

$$-2 \sin(\phi_i(k)) \sin(\phi_j(k)) - 2 \cos(\phi_i(k)) \cos(\phi_j(k))\} \quad (4.11)$$

$$= \frac{1}{2} \sum_{k \in \mathcal{P}} (\cos(\phi_i) - \cos(\phi_j))^2 + (\sin(\phi_i) - \sin(\phi_j))^2 \quad (4.12)$$

$$= \frac{1}{2} \|e^{j\phi_i} - e^{j\phi_j}\|^2 \quad (4.13)$$

To recapitulate, computing the orientation of gradient for an image, and picking its complex exponential as a new representation,

$$\mathbf{z}_i = e^{j\phi_i}, \quad (4.14)$$

provides a set of features for which the typical  $l_2$ -distance is robust with respect to outliers in the original domain, i.e. the domain of pixel values.

Furthermore, PCA can be carried out in the domain of the transformed data  $\mathbf{z}_i$  by looking for a set of orthonormal vectors  $\mathbf{U} = [\mathbf{u}_1 | \dots | \mathbf{u}_K]$  which minimizes

$$\|\mathbf{Z} - \mathbf{U}\mathbf{U}^*\mathbf{Z}\|_F^2, \quad (4.15)$$

where  $\mathbf{U}^*$  denotes the conjugate transpose of  $\mathbf{U}$  and  $\mathbf{Z}$  is the matrix with the available measurements in its columns. This procedure is connected to the version of robust PCA introduced in [30] and yields a subspace of small dimensionality onto which to project the data while retaining most of its “real” variance, and disregarding the variance deriving from outliers. Without going into excessive detail, this desirable property is inherited from the norm minimized in (4.15), computed in the gradient orientation domain, which we have shown to be stable in the presence of corrupted pixel values.

## 4.3 Simulations

We evaluate the proposed model on a set of experiments comprising synthetic and real data.

### 4.3.1 Synthetic data

As a motivational synthetic experiment, we create some artificial data according to model (3.1)-(3.6). The data, 400-dimensional, is folded into 20 x 20 pixel images for an easier visualization. The real values for the factor matrices  $\mathbf{F}$  and  $\mathbf{G}$ , depicted in the first row

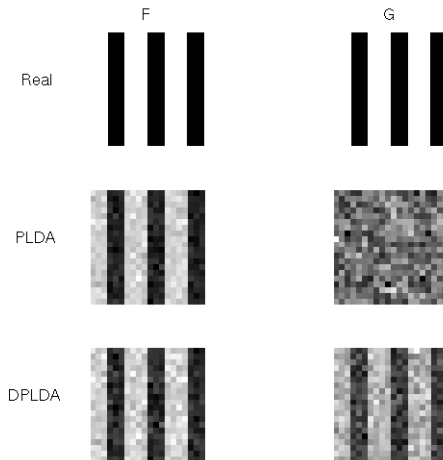


Figure 4.2: Private and shared factor matrices learnt by PLDA and DPLDA compared with the original values.

of figure 4.2, are a series of ones and zeros alternated, so as to create eight strips in the resulting images. Note that in this specific case  $\mathbf{F} = \mathbf{G}$ . For simplicity, the observation noise process is made isotropic with covariance matrix  $\Sigma = \sigma \mathbf{I}$ . The matrix  $\mathbf{A}$ , associated with the linear evolution of the within-individual latent variable, is set to  $\mathbf{A} = 0.98$ , mimicking the smooth variations encountered in real videos.

We then assess the performance of PLDA and DPLDA by learning the necessary parameters for both models, and consequently exploit them for recognition. As argued at the beginning of section 3.1, DPLDA will enforce a proximity constraint on subsequent estimates of the within-individual latent variable  $\mathbf{w}_{ij}^{t-1}$ ,  $\mathbf{w}_{ij}^t$ , due to the dynamical model at its core. PLDA, on the other hand is a static model, and as such it won't take into account any temporal information, so that  $\mathbf{w}_{ij}^t$  is supposed to have zero mean, instead of  $\mathbf{A}\mathbf{w}_{ij}^{t-1}$ . As a result, we expect DPLDA to have a better performance in terms of parameter reconstruction, and consequently of recognition performance.

In our synthetic experiment we observe that, if the observation variance  $\sigma$  is kept low, both models learn matrices  $\mathbf{F}$  and  $\mathbf{G}$  that are very similar to the original versions, consequently performing well in recognition. As the variance increases, however, PLDA's performance worsens quickly, until the estimate for matrix  $\mathbf{G}$  deteriorates completely, as depicted in the second row of figure 4.2. The model thus fails to ascribe the observations to either private or shared latent space in the right proportion, which in turn causes a drop in recognition performance. DPLDA, on the other hand, when dealing with the same extent of noise, still manages to reconstruct successfully both factor matrices and retains a good accuracy, as depicted in the third row of figure 4.2.

### 4.3.2 Real data

#### Verification

As a first experiment with real data, we compare the performance of classical LDA, PLDA and DPLDA in a verification task on the YouTube Faces database [35]. The data consists of 3,425 videos of 1,595 different people, all of which were downloaded from Youtube. An average of 2.15 videos are available per subject (ranging from 1 to 6) with a mean duration

---

Method	Accuracy $\pm$ SE
LDA	$0.723 \pm 0.0054$
PLDA	$0.830 \pm 0.0091$
DPLDA	$0.845 \pm 0.0065$

Table 4.1: YTF verification. Average accuracy and standard error.

of 181 frames.

A benchmark protocol is defined for verification tasks within the database. More precisely, 5000 pairs of videos were selected, half of them depicting the same person, the other half depicting different people. The pairs are further divided into 10 splits, onto which verification has to be performed separately, exploiting the information from the other splits. The restricted protocol only allows access to this information, whereas the unrestricted protocol allows to incorporate information about the identity of the subjects during the training procedure. Since we’re testing class-based methods derived from LDA, we resort to the latter protocol.

We select all the people with 4 or more videos and exploit them for training. However the training set is made split-dependent by erasing from it all videos of all the people that appear in any of the 500 pairs of the split itself. This way we can carry out a more significant test, less prone to overfitting and more indicative of the real performance achieved by each method. Typically the training set for each split consists of roughly 200 different identities, with 4-6 videos each. Again, none of the people depicted in the training set used for a split appear in any video of the respective split.

We consider the first 60 frames for the videos in the database; for each frame we obtain a frontalized version through the FAR algorithm. Consequently, we crop the image from the eyebrows to the lower lip with a bounding box of 65 x 61 pixels. We then perform whitened IGO-PCA on this rectangle of pixel values, reducing dimensionality to 1000 features, which proved to be the lowest size to yield good performance on this database.

As a final step of the verification pipeline we interchange LDA, PLDA and DPLDA, keeping all the other steps identical. For LDA, we compute the distance between all the frames of the first video in a pair from all the frames in the second video of the same pair. Their average is regarded as a video-to-video distance. When considering a specific split, we incorporate the same/not-same information from the remaining 9 splits to learn a linear distance-based classifier, which we then apply to the distances in the current split, determining the final confirm/reject choice.

For PLDA we adopt a similar fusion metrics, applying it to both the sum of the likelihoods of two frames taken separately – one from video 1, one from video 2 – and to the joint likelihood under the hypothesis that they share the same identity. We take their difference as a sufficient statistics for which we learn a confirm/reject threshold, based on the remaining 9 splits, and apply it to the current one.

DPLDA was applied in its most basic version, with  $\mathbf{A}_{ij} = \mathbf{A}$ , i.e. estimating a single state evolution matrix supposed valid for every video. Thresholding is performed analogously to PLDA and no score fusion is necessary since the algorithm acts directly on whole videos instead of single frames. Separate results are computed for each split, and finally the 10 verification rates are averaged. We repeated this experiment 5 times, and the mean values are reported in table 4.1.

Note that both DPLDA and PLDA perform significantly better than simple LDA. Also, despite the simplification  $\mathbf{A}_{ij} = \mathbf{A}$ , DPLDA outperforms PLDA by more than 1%. This is

a first important confirmation that explicitly modelling the dynamics for the private latent variable in video sequences, even with a simple linear evolution dependent on a matrix shared among all the videos, can improve the discriminatory power of PLDA models in a significant and challenging setting.

### Identification

We assembled an “in-the-wild” database consisting of 250 videos of 50 people, each depicted in exactly 5 videos. The videos, featuring famous people not present in the YTF database, were downloaded from Youtube at 24fps in medium quality. For each video we manually checked the identity of the person depicted and the quality of the whole sequence, in order to guarantee the presence of significant temporal information to exploit, and avoid still-image slide-shows. Exploiting the pipeline described above, we test PLDA and DPLDA for closed-set and open-set identification on this database.

5 analogous closed-set identification tasks are carried out by picking one of the 5 videos for each person to incorporate in the probe and utilizing the remaining 4 for training. Fixing the size of the shared subspace to 50, we let the size of the private subspace vary from 1 to 50. The 5 tests are repeated 3 times each, with the average results depicted in figure 4.3. For  $\text{Dim}(\mathbf{G}) < 10$ , PLDA performs better than DPLDA, however the latter exhibits a steeper performance increase as  $\text{Dim}(\mathbf{G})$  becomes larger. For  $\text{Dim}(\mathbf{G}) > 15$ , DPLDA consistently outperforms PLDA by roughly 2%.

For open set recognition, we keep the same setting, simply adding a variable number of external videos to the probe. For convenience we picked them from the YFT database, with which our database has no overlap whatsoever. As a performance measurement we pick the generalized identification accuracy, by simply contemplating an additional external class for which no training data is available. We test for different levels of “openness” by letting the percentage of external videos in the probe vary from 17 % to 66 %. For every probe video, we compute the conditional likelihood of it sharing its identity with all the individuals in the gallery, storing the maximum value, along with the marginal likelihood of the probe video alone. The probe video is labelled as external if the difference between conditional and marginal likelihoods is smaller than a specific threshold, otherwise it is given the label that maximizes the conditional likelihood. No specific threshold is learnt, rather we depict the accuracy rate as a function of the threshold values for both PLDA and DPLDA in figure 4.3. Note that the  $x$  axis scale has been normalized so that its extrema actually represent the minimum and maximum values measured for the difference between conditional and marginal likelihood on any probe video. In other words, the graphs are scaled and aligned so that outside the depicted threshold values they are completely flat.

When the threshold is at a minimum, all videos are labelled as internal, hence the advantage of DPLDA due to better performance in closed-set identification. As the threshold increases, more and more videos are rejected as external, until all of them are when the normalized threshold is equal to 1. To this extremum, PLDA and DPLDA have an identical performance in terms of accuracy, which is equal to the percentage of external videos in the probe. For intermediate threshold values, the advantage of DPLDA over PLDA in terms of accuracy is quite evident, thus confirming again the suitability of the proposed method for different task of identity inference in video sequences.

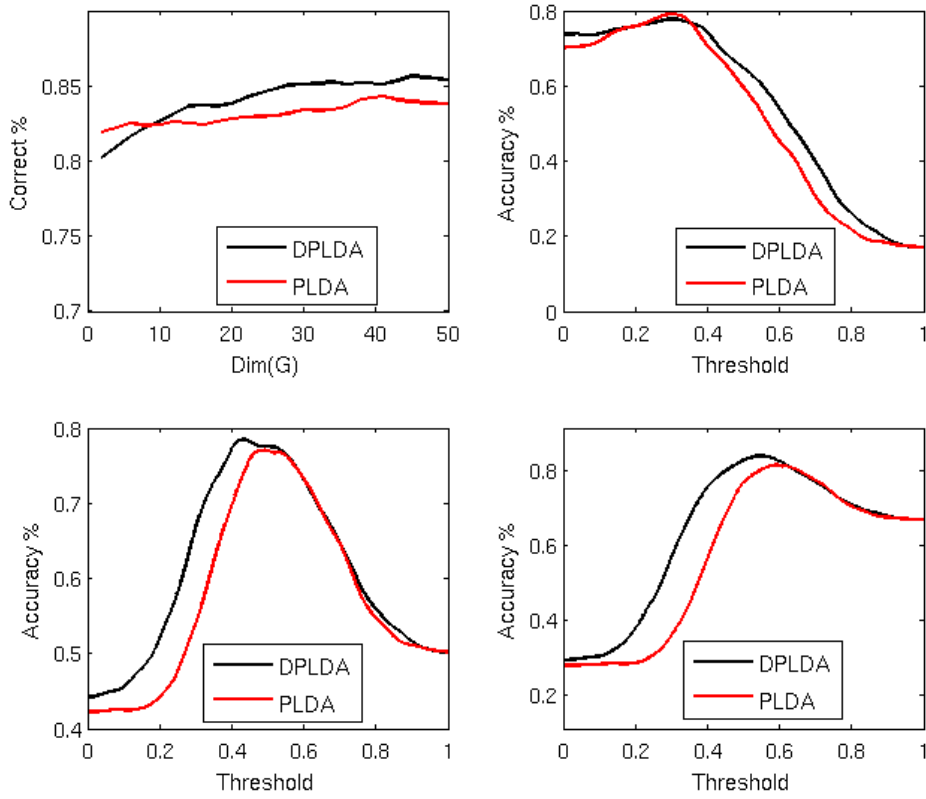


Figure 4.3: Closed- and Open-set identification on our database. (Top-left) Comparison of Closed-set identification accuracy as the private latent space dimensionality increases. (Top-right), (Bottom-left), (Bottom-right) Comparison of Open-set identification accuracy dependent on threshold value. The percentage of external videos are 17%, 50% and 66% respectively.

## 5. Conclusion

An increasing interest has been devoted to face recognition in the last few decades, starting from tasks on still images, progressively generalized to video sequences. Despite the maturity of this technology, there is still a variety of applications susceptible to improvement. Challenges arise especially in the context of non-cooperative video recording, such as ubiquitous CCTV cameras, whereby one or more factors such as quality, pose, illumination, occlusions and duration are uncontrolled.

In this document we focused our attention on in-the-wild databases with uncontrolled acquisition. We introduced a 3-stage pipeline featuring frontalization, feature extraction and recognition to handle this challenging setting. For the recognition stage we chose PLDA for its inherent flexibility, which allowed us to handle the classical tasks of verification and closed-set identification, but also clustering and open-set identification. The proposed recognition pipeline achieved a good performance on the YTF database and on a privately assembled one, only exploiting the information in the first 3 seconds of each video.

Furthermore we leveraged a sequence-based, spatio-temporal approach to propose the first dynamical version of PLDA. We outlined its theoretical advantage over classical PLDA and proved it in three tasks on two different databases, both featuring challenging conditions.

The systematic advantage of DPLDA emerged despite the simple setting we adopted, forcing all videos to share the same matrix for state dynamics ( $\mathbf{A}_{ij} = \mathbf{A}$ ). A richer model, featuring different evolution matrices, is likely to improve performance even further, allowing us to take into account the peculiar dynamics of each video. However, the inversion of the innovation covariance matrices from the Kalman filters represents a bottleneck, hindering the adoption of more complex, performance-improving models. An interesting idea would be to explore approximate approaches to Kalman filtering, such as the conjugate gradient, to circumvent this issue.





# A. The EM Algorithm

The EM algorithm [8] is a powerful statistical tool to perform maximum-likelihood estimates for parameters in models comprising latent variables. It was formalized in a classical paper by Dempster et al. in 1977. The authors recognized similar procedures had been proposed before, as solutions to specific statistical models, which, however, lacked a unifying framework. Notable examples are the Baum-Welch algorithm and K-means clustering.

In this appendix we complement the intuitive explanation provided in section 2.4 with a formal proof of convergence. More specifically, when looking for maximum-likelihood solutions, the aim is to maximize the likelihood of the observations, or equivalently its logarithm

$$\ln P(\mathbf{X}|\boldsymbol{\theta}) = \ln \left[ \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \right], \quad (\text{A.1})$$

where  $\mathbf{Z}$  groups all the latent variables in the model. For any PDF  $q(\mathbf{Z})$  defined over the latent variables the following decomposition holds:

$$\ln P(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + KL(q||p). \quad (\text{A.2})$$

Here  $KL(q||p)$  is the Kullback-Leibner divergence, defined as

$$KL(q||p) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} \right\}, \quad (\text{A.3})$$

which measures the difference between  $q(\mathbf{Z})$  and  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ . As a metric for the distance between PDFs, it is always non-negative, equal to zero if and only if  $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ .

$\mathcal{L}(q, \boldsymbol{\theta})$  is a functional of  $q(\mathbf{Z})$  and a function of  $\boldsymbol{\theta}$  defined as

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right\}, \quad (\text{A.4})$$

representing a lower bound on the log-likelihood. As we will show, iterative optimization of  $\mathcal{L}(q, \boldsymbol{\theta})$  results in the maximization of the log-likelihood, hence in a ML solution.

In the E-step we optimize  $\mathcal{L}(q, \boldsymbol{\theta}^{old})$  with respect to  $q(\mathbf{Z})$  while holding parameter  $\boldsymbol{\theta}$  fixed to a previous estimate  $\boldsymbol{\theta}^{old}$ . Since  $\ln p(\mathbf{X}|\boldsymbol{\theta})$  does not depend on  $q(\mathbf{Z})$ , the largest values of  $\mathcal{L}(q, \boldsymbol{\theta}^{old})$  will occur when the Kullback-Leibner divergence is at its smallest, which happens for  $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})$ .

During the M-step,  $q(\mathbf{Z})$  is held fixed and  $\mathcal{L}(q, \boldsymbol{\theta})$  is maximized with respect to  $\boldsymbol{\theta}$ . Unless  $\mathcal{L}(q, \boldsymbol{\theta})$  has a maximum in  $\boldsymbol{\theta}^{old}$ , the new optimum  $\boldsymbol{\theta}^{new}$  will result in an increase of  $\mathcal{L}(q, \boldsymbol{\theta})$  and an even larger increase of the log-likelihood, due to the fact that  $KL(q||p) = KL(p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{new}))$  will become positive.

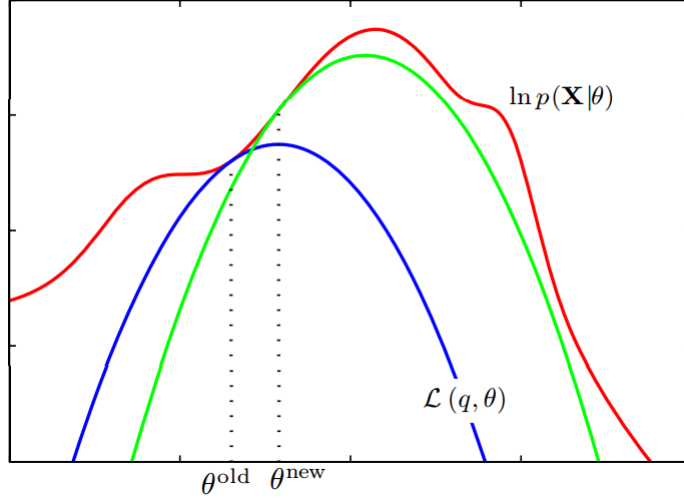


Figure A.1: Graphical depiction of the iterative optimization carried out through the EM algorithm.

Furthermore, note that the following decomposition holds

$$\begin{aligned}\mathcal{L}(q, \boldsymbol{\theta}) &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) \\ &= Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) + \text{const},\end{aligned}\tag{A.5}$$

where the constant is the entropy of the PDF  $q(\mathbf{Z})$  and is independent of  $\boldsymbol{\theta}$ . To summarize, maximization (2.23) is equivalent to optimizing  $\mathcal{L}(q, \boldsymbol{\theta})$ , which, in alternation with the evaluation of  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ , guarantees a continuous increase in the likelihood of the observations until convergence to a local maximum.

This process is depicted in figure A.1, taken from [6] where we suppose to start from an initial parameter value  $\boldsymbol{\theta}^{old}$ . The curve in red represents the log-likelihood to be maximized, the curves in blue and green is the auxiliary functional  $\mathcal{L}(q, \boldsymbol{\theta})$  we optimize instead. In the first E-step, we compute  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ , giving rise to  $\mathcal{L}(q, \boldsymbol{\theta})$  which equals the log-likelihood at  $\boldsymbol{\theta}^{old}$ . The lower bound  $\mathcal{L}(q, \boldsymbol{\theta})$  is then maximized with respect to  $\boldsymbol{\theta}$  during the M-step, which gives a new parameter value  $\boldsymbol{\theta}^{new}$ , and yields an improvement in the log-likelihood. In the subsequent E-step,  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$  is evaluated for the new set of parameters  $\boldsymbol{\theta}^{new}$ , again providing a convex lower bound represented by the green curve.

# Acknowledgments

I am pleased to thank my official supervisor from Imperial College, Stefanos Zafeirou, and my unofficial supervisor, Mihalis Nicolau, for their patience, their help, their guidance.

Nice things should not be left unspoken, the following should be uttered in Italian.

Ringrazio il mio relatore patavino, Ruggero Carli, per la sua disponibilità e i suoi preziosi consigli.

Un sentito ringraziamento va ai miei amici e conterranei Giorgia, Irene, Ama, l'Ale, Prince, Gaspa, Chicco, Checco, Cant. Grazie per esserci stati, perché ci siete e ci sarete. Grazie perché rendete casa un posto meraviglioso a cui fare ritorno.

Un grazie ai *Sattoni* perché cambiano un po' ogni anno, ma in fondo non cambiano mai. Grazie per le risate, le pallavolate, per la quiete e le monete. È vero.

Ringrazio Paola, che mi ha cresciuto, educato, amato. Grazie per avermi trattato come un figlio con discrezione ed affetto ineguagliabili.

Un grazie immenso va a mio padre, Ennio, mio primo sostenitore da sempre. Grazie per tutto quello che hai fatto, grazie per gli abbracci, per i consigli, per le sberle meritate. Grazie per aver fatto del tuo meglio per insegnarmi a stare al mondo; il risultato, colui che scrive, spero non sia così male.

Da ultimo ringrazio mia sorella, Olga, che mi ha preso per mano nel giorno della mia nascita e da allora non ha mai lasciato andare. Grazie perché sei un modello, una madre, una sorellina, un'amica. Grazie perché in me c'è così tanto di te, che la mia vita senza te proprio non riesco a immaginarla.



# Bibliography

- [1] M. Al-Azzeh, A. Eleyan, and H. Demirel. Pca-based face recognition from video using super-resolution. *Proc. 2008 International Symposium on Computer and Information Sciences*, pages 1 – 4, 2008.
- [2] O. Arandjelović and R. Cipolla. Probabilistic models for inference about identity. *Image and Vision Computing*, 24:639 – 647, 2006.
- [3] J. R. Barr, K. W. B. P. J. FLynn, and S. Biswas. Face recognition from video: a review. *International Journal of Pattern Recognition and Artificial Intelligence*, 2012.
- [4] L. Benedikt, D. Cosker, P. Rosin, and D. Marshall. 3d facial gestures in biometrics: from feasibility study to application. *Proc. 2008 IEEE International Conference on Biometrics: Theory, Applications and Systems*, 34:1 – 6, 2008.
- [5] S. Berrani and C. Garcia. Enhancing face recognition from video sequences using robust statistics. *Proc. 2005 IEEE Conference on Advanced Video and Signal Based Surveillance*, 34:324 — 329.
- [6] C. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer, 2006.
- [7] L. Chen, H. Liao, and J. Lin. Person identification using facial motion. *Proc. 2001 International Conference on Image Processing*, 2:677 – 680, 2001.
- [8] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1 – 38, 1977.
- [9] R. A. Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7(2):179 – 188, 1936. ISSN 2050-1439. doi: 10.1111/j.1469-1809.1936.tb02137.x. URL <http://dx.doi.org/10.1111/j.1469-1809.1936.tb02137.x>.
- [10] A. C. G. Aggarwal and R. Chellappa. A system identification approach for video-based face recognition. *Proc. 2004 International Conference on Pattern Recognition*, 4:175 – 178, 2004.
- [11] D. Gorodnichy. Video-based framework for face recognition in video. *Proc. 2005 Canadian Conference on Computer and Robot Vision*, 34:330 – 338, 2005.
- [12] B. Gunturk, A. Batur, Y. Altunbasak, M. Hayes, and R. Mersereau. Eigenfacedomain super-resolution for face recognition. *IEEE Transactions on Image Processing*, pages 587 – 606, 2003.

- 
- [13] A. Hadid and M. Pietikäinen. Selecting models from videos for appearance-based face recognition. *Proc. 17th International Conference on Pattern Recognition*, 1:304 – 308, 2004.
- [14] A. Hadid and M. Pietikäinen. Combining appearance and motion for face and gender recognition from videos. *IEEE Transactions on Pattern Recognition and Machine Intelligence*, 42:2818 – 2827, 2009.
- [15] R. R. Jillela and A. Ross. Adaptive frame selection for improved face recognition in low-resolution videos. *Proc. 2009 International Joint Conference on Neural Networks*, pages 2835 – 2841, 2009.
- [16] P. Li, Y. Fu, U. Mohammed, J. H. Elder, and S. J. Price. Probabilistic models for inference about identity. *IEEE Transactions on Pattern Recognition and Machine Intelligence*, 34:144 – 157, 2012.
- [17] Y. Li, S. Gong, and H. Liddell. Constructing facial identity surfaces for recognition. *International Journal of Computer Vision*, 53:71 – 92, 2003.
- [18] Lina, T. Takahashi, I. Ide, and H. Murase. Incremental unsupervised-learning of appearance manifold with view-dependent covariance matrix for face recognition from video sequences. *IEICE Transactions on Information and Systems*, 34:642 – 652, 2009.
- [19] X. Liu and T. Chen. Face mosaicing for pose robust video-based recognition. *Proc. of the 8th Asian Conference on Computer Vision*, pages 662 – 671, 2007.
- [20] X. Liu and T. Cheng. Video-based face recognition using adaptive hidden markov models. *Proc. 2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 1:340 – 345, 2003.
- [21] F. Matta and J.-L. Dugelay. Video face recognition: A physiological and behavioural multimodal approach. *Proc. 2007 IEEE International Conference on Image Processing*, pages 497 – 500, 2007.
- [22] U. Park, H. Chen, and A. Jain. 3d model-assisted face recognition in video. *Proc. 2005 Canadian Conference on Computer and Robot Vision*, pages 322 – 329, 2005.
- [23] A. Ramanathan, R. Chellappa, and S. Biswas. Computational methods for modeling facial aging: A survey. *Journal of Visual Languages and Computing*, 20:131 – 144, 2009.
- [24] M. Ribeiro. Kalman and extended Kalman filters: Concept, derivation and properties. Technical report., Institute for Systems and Robotics, Lisboa, 2004.
- [25] C. Sagonas, Y. Panagakis, S. Zafeiriou, and M. Pantic. Face frontalization for alignment and recognition. (to appear).
- [26] L. Sirovich and M. Kirby. Low-dimensional procedure for characterization of human faces. *Journal of Optical Society of America*, 3:519 – 524, 1987.
- [27] J. Stallkamp, H. K. Ekenel, and R. Stiefelhagen. Video-based face recognition on real-world data. *Proc. 2007 IEEE International Conference on Computer Vision*, 206:1 – 8.

- [28] D. Thomas, K. W. Bowyer, and P. J. Flynn. Multi-frame approaches to improve face recognition. *Proc. 2007 IEEE Workshop on Motion and Video Computing*, page 19, 2007.
- [29] M. E. Tipping and C. M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society, Series B*, 61:611 – 622, 1999.
- [30] F. D. L. Torre and M. Black. A framework for robust subspace learning. *IJCV* 54, pages 117 – 142, 2003.
- [31] G. Tzimiropoulos, S. Zafeiriou, and M. Pantic. Subspace learning from image gradient orientations. *IEEE TPAMI*, (to appear).
- [32] R. Wang, S. Shan, X. Chen, and W. Gao. Manifold-manifold distance with application to face recognition based on image set. *Proc. 2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1 – 8, 2008.
- [33] M. E. Wibowo, D. Tjondronegoro, and V. Chandran. Probabilistic matching of image sets for video-based face recognition. *Digital Image Computing Techniques and Applications (DICTA)*, pages 1 – 6, 2012.
- [34] M. E. Wibowo, D. Tjondronegoro, L. Zhang, and I. Himawan. Heteroscedastic probabilistic linear discriminant analysis for manifold learning in video-based face recognition. *Applications of Computer Vision*, pages 46 – 52, 2013.
- [35] L. Wolf, T. Hassner, and I. Maoz. Face recognition in unconstrained videos with matched background similarity. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2011.
- [36] Q. Xiong and C. Jaynes. Mugshot database acquisition in video surveillance networks using incremental auto-clustering quality measures. *Proc. 2003 IEEE Conference on Advanced Video and Signal Based Surveillance*, page 191.
- [37] O. Yamaguchi, K. Fukui, and K. Maeda. Face recognition using temporal image sequence. *Proc. 1998 IEEE International Conference on Automatic Face and Gesture Recognition*, pages 318 – 323, 1998.
- [38] Y. Zhang and A. M. Martinez. A weighted probabilistic approach to face recognition from multiple images and video sequences. *Image and Vision Computing*, 24:626 — 638, 2006.
- [39] X. Zhou and B. Bhanu. Human recognition based on face profiles in video. *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, page 15, 2005.