

Simultaneous and Causal Appearance Learning and Tracking

J. Melenchón, I. Iriondo and L. Meler

*Communications and Signal Theory Department,
Enginyeria La Salle, Universitat Ramon Llull, Pg. Bonanova 8, 08022, Barcelona, Spain*

Received 20 December 2004 ; accepted 22 March 2005

Abstract

A novel way to learn and track simultaneously the appearance of a previously non-seen face without intrusive techniques can be found in this article. The presented approach has a causal behaviour: no future frames are needed to process the current ones. The model used in the tracking process is refined with each input frame thanks to a new algorithm for the simultaneous and incremental computation of the singular value decomposition (SVD) and the mean of the data. Previously developed methods about iterative computation of SVD are taken into account and an original way to extract the mean information from the reduced SVD of a matrix is also considered. Furthermore, the results are produced with linear computational cost and sublinear memory requirements with respect to the size of the data. Finally, experimental results are included, showing the tracking performance and some comparisons between the batch and our incremental computation of the SVD with mean information.

1 Introduction

The last years have witnessed extraordinary advances in computer and communications technology, leading to an increasing availability of information and processing capabilities of multimedia data [1], [2]. This fact is resulting in a higher and wider demand for easier access to information [3]. On one hand, this information is mainly stored in digital format, so its access is limited to the user's ability to communicate with computers. On the other hand, it has been remarked the great expressive power of the natural language used in human-human communication, as well as its intrinsic multimodal features [4]. Consequently, the access to digital information could be carried out using this natural language: reducing the necessity of knowing a specific way to interact with the computer and taking advantage of its expressive features. Moreover, multimodal interfaces with an audio visual system like a talking head could be used in order to speak to the user in natural language. As a result, talking heads used in multimodal interfaces seem to be a proper solution for making access to information easier and more pleasing for human users.

As explained in [4], multimodal input analysis is necessary when working with multimodal interfaces and relies on interaction devices e.g. facial trackers. Some non-intrusive visual trackers can be used in this scheme because they retain information regarding to position, scale, orientation and appearance of the tracked element, e.g. [5], [6], [7], [8] and [9]. Nevertheless, the whole sequence is needed by these algorithms to be processed off-line (they have a non-causal behaviour); as a result, a real time implementation of these methods is impossible, even without considering their computational cost. This temporal restriction is caused by the computation

Correspondence to: jmelen@salleurl.edu

Recommended for acceptance by Perales F., Drapper B.

ELCVIA ISSN:1577-5097

Published by Computer Vision Center / Universitat Autònoma de Barcelona, Barcelona, Spain

of a Singular Value Decomposition (SVD) over the whole observed data. Moreover, memory resources are greatly affected by this fact, limiting the duration of the observed sequence. Incremental SVD computation techniques as [10], [11] and [12] may be useful in this case, but they do not take into consideration the mean of the data, which is crucial in the classification of the different gestures. Fortunately, this is taken into account in [13] and [14]. By one hand, the work presented in [13] does not propose a method to extract the mean information from a given SVD and it can only update the SVD from two other known SVD. By the other hand, Skočaj presented in [14] a method with a similar performance to the one achieved in this paper, but he focused on incremental Principal Component Analysis rather than incremental SVD.

In this paper, a new method for updating both SVD and mean information as well as extracting the mean of the data contained in a given SVD without increasing the cost order of either time or memory is presented in Sect. 2. The application of this new method is carried out in Sect. 3 by a causal algorithm for the tracking and learning of the facial appearance of a person. Experimental results are given in Sect. 4 and concluding remarks are explained in Sect. 5.

2 Incremental SVD with Mean Update

2.1 Fundamentals

The singular value decomposition of matrix $\mathbf{M}_{p \times q} = [\mathbf{m}_1 \cdots \mathbf{m}_q]$ is given by:

$$\mathbf{M}_{p \times q} = \mathbf{U}_{p \times p} \mathbf{\Sigma}_{p \times q} \mathbf{V}_{q \times q}^T, \quad (1)$$

where $\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_p]$ and $\mathbf{V} = [\mathbf{v}_1 \cdots \mathbf{v}_q]$ are orthonormal matrices; \mathbf{u}_i are the eigenvectors of $\mathbf{M}\mathbf{M}^T$ and span the column space of \mathbf{M} ; \mathbf{v}_i are the eigenvectors of $\mathbf{M}^T\mathbf{M}$ and span the row space of \mathbf{M} ; and $\mathbf{\Sigma}$ is a diagonal matrix with the singular values of either $\mathbf{M}\mathbf{M}^T$ and $\mathbf{M}^T\mathbf{M}$ in descending order. Notice that if \mathbf{M} is a rank r matrix, where $r \leq p$ and $r \leq q$, its corresponding $\mathbf{\Sigma}$ has only r non-null singular values and (1) can be rewritten as the *thin SVD*: $\mathbf{M}_{p \times q} = \mathbf{U}_{p \times r} \mathbf{\Sigma}_{r \times r} \mathbf{V}_{q \times r}^T$. By the other hand, let $\mathbf{C}_{r \times q} = \mathbf{U}_{p \times r}^T \mathbf{M}_{p \times q}$ be the projections of the columns of \mathbf{M} over the eigenspace spanned by \mathbf{U} . Using the *thin SVD* expression the projections matrix $\mathbf{C} = [\mathbf{c}_1 \cdots \mathbf{c}_q]$ can be written also as $\mathbf{C}_{r \times q} = \mathbf{\Sigma}_{r \times r} \mathbf{V}_{q \times r}^T$.

In other fields, like classification problems pointed by [13], a more suitable representation of \mathbf{M} can be achieved including mean information $\overline{\mathbf{m}} = \frac{1}{q} \sum_{i=1}^q \mathbf{m}_i$ in (1), which has to be computed and subtracted previously from \mathbf{M} in order to be able to generate the SVD of $\mathbf{M} - \overline{\mathbf{m}} \cdot \mathbf{1}$:

$$\mathbf{M}_{p \times q} = \mathbf{U}_{p \times r} \mathbf{\Sigma}_{r \times r} \mathbf{V}_{q \times r}^T + \overline{\mathbf{m}}_{p \times 1} \mathbf{1}_{1 \times q}. \quad (2)$$

2.2 Updating SVD

Assuming an existing SVD (1), if new columns $\mathbf{I}_{p \times c} = [\mathbf{I}_1 \cdots \mathbf{I}_c]$ are added in order to obtain a new matrix $\mathbf{M}'_{p \times (q+c)} = [\mathbf{M}_{p \times q} \quad \mathbf{I}_{p \times c}]$, the SVD of \mathbf{M}' can be updated from (1) using methods like [11] and [12], achieving:

$$\mathbf{M}'_{p \times (q+c)} = \mathbf{U}'_{p \times r'} \mathbf{\Sigma}'_{r' \times r'} \mathbf{V}'_{(q+c) \times r'}{}^T. \quad (3)$$

Otherwise, if the representation of \mathbf{M}' is chosen to be as (2) and $\overline{\mathbf{m}}'$ is set to $\frac{1}{q+c} (\sum_{k=1}^q \mathbf{m}_k + \sum_{l=1}^c \mathbf{I}_l)$ the SVD becomes:

$$\mathbf{M}'_{p \times (q+c)} = \mathbf{U}'_{p \times r'} \mathbf{\Sigma}'_{r' \times r'} \mathbf{V}'_{(q+c) \times r'}{}^T + \overline{\mathbf{m}}'_{p \times 1} \mathbf{1}_{1 \times (q+c)}. \quad (4)$$

Starting from (2) and matrix \mathbf{I} , (4) can be obtained using the method proposed by [13] if the SVD of \mathbf{I} is previously computed and q and c are known beforehand. A new method for updating both the SVD and the mean using only the new observations and previous factorization is presented in Sect. 2.3.

2.3 Updating SVD and Mean

Beginning with an existing factorization of \mathbf{M}_i as in (5), it is desired to obtain the SVD and mean of \mathbf{M}_f shown in (6):

$$\mathbf{M}_i = \mathbf{U}_i \boldsymbol{\Sigma}_i \mathbf{V}_i^T + \bar{\mathbf{m}}_i \mathbf{1} . \quad (5)$$

$$\mathbf{M}_f = [\mathbf{M}_i \quad \mathbf{I}] = \mathbf{U}_f \boldsymbol{\Sigma}_f \mathbf{V}_f^T + \bar{\mathbf{m}}_f \mathbf{1} . \quad (6)$$

Defining $\hat{\mathbf{M}}_i$ (7) and centering new columns \mathbf{I} around $\bar{\mathbf{m}}_i$ (8), it can be written:

$$\hat{\mathbf{M}}_i = \mathbf{M}_i - \bar{\mathbf{m}}_i \mathbf{1} = \mathbf{U}_i \boldsymbol{\Sigma}_i \mathbf{V}_i^T . \quad (7)$$

$$[\mathbf{M}_i \quad \mathbf{I}] - \bar{\mathbf{m}}_i \mathbf{1} = \mathbf{U}_f \boldsymbol{\Sigma}_f \mathbf{V}_f^T + \bar{\mathbf{m}}_f \mathbf{1} - \bar{\mathbf{m}}_i \mathbf{1} . \quad (8)$$

$$[\mathbf{M}_i - \bar{\mathbf{m}}_i \mathbf{1} \quad \mathbf{I} - \bar{\mathbf{m}}_i \mathbf{1}] = \mathbf{U}_f \boldsymbol{\Sigma}_f \mathbf{V}_f^T + (\bar{\mathbf{m}}_f - \bar{\mathbf{m}}_i) \mathbf{1} . \quad (9)$$

$$[\hat{\mathbf{M}}_i \quad \hat{\mathbf{I}}] = \mathbf{U}_t \boldsymbol{\Sigma}_t \mathbf{V}_t^T . \quad (10)$$

The new columns $\mathbf{I}_{p \times c}$ (see sect. 2.2) will be known through this paper as the *update block*. Note that (10) is the updated SVD from (7) when some new observations $\hat{\mathbf{I}}$ are added. This update can be done as [12] suggests:

$$[\hat{\mathbf{M}}_i \quad \hat{\mathbf{I}}] = [\mathbf{U}_i \quad \mathbf{Q}_i] \begin{bmatrix} \boldsymbol{\Sigma}_i & \mathbf{U}_i^T \hat{\mathbf{I}} \\ \mathbf{0} & \mathbf{Q}_i^T \hat{\mathbf{I}} \end{bmatrix} \begin{bmatrix} \mathbf{V}_i^T & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} = [\mathbf{U}_i \quad \mathbf{Q}_i] \mathbf{U}_d \boldsymbol{\Sigma}_d \mathbf{V}_d^T \begin{bmatrix} \mathbf{V}_i^T & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} = \mathbf{U}_t \boldsymbol{\Sigma}_t \mathbf{V}_t^T \quad (11)$$

where QR-decomposition is done to $\hat{\mathbf{I}} - \mathbf{U}_i \mathbf{U}_i^T \hat{\mathbf{I}} = \mathbf{Q}_i \mathbf{R}_i$ to obtain an orthogonal basis \mathbf{Q}_i for the reconstruction error. Next, the *mean update algorithm* can be executed starting from the knowledge of $\hat{\mathbf{V}}_t^T = \hat{\mathbf{V}}_t^T + \bar{\mathbf{v}}_t \mathbf{1}$, where $\bar{\mathbf{v}}_t = \frac{1}{q+c} \sum_{k=1}^{q+c} \mathbf{v}_k$:

$$[\hat{\mathbf{M}}_i \quad \hat{\mathbf{I}}] = \mathbf{U}_t \boldsymbol{\Sigma}_t \hat{\mathbf{V}}_t^T + \mathbf{U}_t \boldsymbol{\Sigma}_t \bar{\mathbf{v}}_t \mathbf{1} = \mathbf{U}_t \boldsymbol{\Sigma}_t \hat{\mathbf{V}}_t^T + \bar{\mathbf{m}}_t \mathbf{1} . \quad (12)$$

$$[\hat{\mathbf{M}}_i \quad \hat{\mathbf{I}}] = \mathbf{U}_t \boldsymbol{\Sigma}_t \mathbf{R}_v^T \mathbf{Q}_v^T + \bar{\mathbf{m}}_t \mathbf{1} = \mathbf{U}_f \boldsymbol{\Sigma}_f \mathbf{V}_u^T \mathbf{Q}_v^T + \bar{\mathbf{m}}_t \mathbf{1} = \mathbf{U}_f \boldsymbol{\Sigma}_f \mathbf{V}_f^T + \bar{\mathbf{m}}_t \mathbf{1} . \quad (13)$$

$$[\hat{\mathbf{M}}_i \quad \hat{\mathbf{I}}] + \bar{\mathbf{m}}_i \mathbf{1} = \mathbf{U}_f \boldsymbol{\Sigma}_f \mathbf{V}_f^T + \bar{\mathbf{m}}_t \mathbf{1} + \bar{\mathbf{m}}_i \mathbf{1} . \quad (14)$$

$$[\mathbf{M}_i \quad \mathbf{I}] = \mathbf{U}_f \boldsymbol{\Sigma}_f \mathbf{V}_f^T + \bar{\mathbf{m}}_f \mathbf{1} . \quad (15)$$

It is assumed $\mathbf{Q}_v \mathbf{R}_v$ as the QR-decomposition of $\hat{\mathbf{V}}_t$, $\mathbf{U}_f \boldsymbol{\Sigma}_f \mathbf{V}_u^T$ as the SVD of $\mathbf{U}_t \boldsymbol{\Sigma}_t \mathbf{R}_v^T$ and $\bar{\mathbf{m}}_f = \bar{\mathbf{m}}_t + \bar{\mathbf{m}}_i$. Note that (15) and (6) are the same expression.

2.4 Mean Extraction from a Given SVD

The previous method can also be used to extract the mean information from an existing SVD, e.g. trying to express $\mathbf{S} = \mathbf{U}_t \boldsymbol{\Sigma}_t \mathbf{V}_t^T$ as $\mathbf{S} = \mathbf{U}_f \boldsymbol{\Sigma}_f \mathbf{V}_f^T + \bar{\mathbf{s}} \cdot \mathbf{1}$ setting $[\hat{\mathbf{M}}_i \quad \hat{\mathbf{I}}] = \mathbf{S}$ and $\bar{\mathbf{m}}_t = \mathbf{0}$ in (12) to (15).

2.5 Time and Memory Complexity

The mean update presented in section 2.3 does not increase the order of resources required in methods of incremental SVD developed in [10], [12], [11], [13] and [14]. The computational cost becomes $O(qr^2 + pr^2)$ and the memory complexity is $O(pr + qr)$, as shown in Table 1.

3 On-the-fly Face Training

In this paper, *On-the-fly Face Training* is defined as the process of learning the photo-realistic facial appearance model of a person observed in a sequence in a rigorous causal fashion. This fact means that it is not necessary to take into account subsequent images when adding the information of the current one, which is considered only once. Note that the facial appearance is learnt in the same order as the captured images, allowing a real-time learning capability in near future, as computational resources are constantly being increased.

Table 1: Resource order requirements of the proposed *mean update algorithm*

Operation	Comp. cost	Mem. requirements
$\mathbf{V}_{q \times r}^T - \left(\frac{1}{q} \sum_{k=1}^q (\mathbf{v}_k)_{r \times 1} \right) \mathbf{1}_{1 \times q} \rightarrow \hat{\mathbf{V}}_{q \times r}^T$	$O(qr)$	$O(qr + r)$
$\hat{\mathbf{V}}_{q \times r} \rightarrow (\mathbf{Q}_v)_{q \times r} (\mathbf{R}_v)_{r \times r}$	$O(qr^2)$	$O(qr + r^2)$
$(\mathbf{U}_i)_{p \times r} (\mathbf{\Sigma}_i)_{r \times r} (\mathbf{R}_v^T)_{r \times r} \rightarrow \mathbf{T}_{p \times r}$	$O(pr^2 + r^3)$	$O(pr + r^2)$
$\mathbf{T}_{p \times r} \rightarrow (\mathbf{U}_f)_{p \times r} (\mathbf{\Sigma}_f)_{r \times r} (\mathbf{V}_u^T)_{r \times r}$	$O(pr^2)$	$O(pr + r^2)$
$(\mathbf{V}_f)_{q \times r} \rightarrow (\mathbf{Q}_v)_{q \times r} (\mathbf{V}_u)_{r \times r}$	$O(qr^2)$	$O(qr + r^2)$
Totals, assuming $p \gg r$ and $g \gg r$	$O(qr^2 + pr^2)$	$O(pr + qr)$

3.1 Data Representation

An N image sequence $\mathbf{S} = [\mathbf{I}_1 \cdots \mathbf{I}_N]$ and a set of four masks $\mathbf{\Pi} = \{\pi^1, \dots, \pi^4\}$, attached to four facial elements (like mouth, eyes or forehead), are given. For each image \mathbf{I}_t , its specific mouth, eyes and forehead appearance are extracted using $\mathbf{\Pi}$, obtaining four observation vectors \mathbf{q}_t^r (see Fig. 1). Therefore, four observation matrices \mathbf{O}^r can be obtained from the application of the set of masks $\mathbf{\Pi}$ over the sequence \mathbf{S} . Dimensionality reduction of \mathbf{O}^r can be achieved using SVD [15]: $\mathbf{O}^r = [\mathbf{o}_1^r \cdots \mathbf{o}_N^r] = \mathbf{U}^r \mathbf{\Sigma}^r (\mathbf{V}^r)^T + \bar{\mathbf{o}}^r \mathbf{1}_{1 \times N}$, where $\bar{\mathbf{o}}^r = \frac{1}{N} \sum_{k=1}^N \mathbf{o}_k^r$. Note that facial element appearances can be parameterized as $\mathbf{C} = \mathbf{\Sigma}^r (\mathbf{V}^r)^T$ (see Sect. 2.1). In the example proposed in this paper, faces composed of 41205 pixels could be codified with 35 coefficients, representing a reduction of more than 99.9% without any loss of perceptual quality (see Fig. 2).

3.2 Training Process

One major drawback of the parametrization presented in section 3.1 consists in the image alignment of the sequence [5]. Unless all face images through the whole sequence have the same position, ghosly results may appear and suboptimal dimensionality reduction will be achieved. The tracking scheme presented in this paper combines simultaneously both processes of learning and alignment. First of all, the four masks π^r are manually extracted from the first image \mathbf{I}_1 of sequence \mathbf{S} and the first observation vectors $\mathbf{o}_1^1, \dots, \mathbf{o}_1^4$ are obtained. Next,

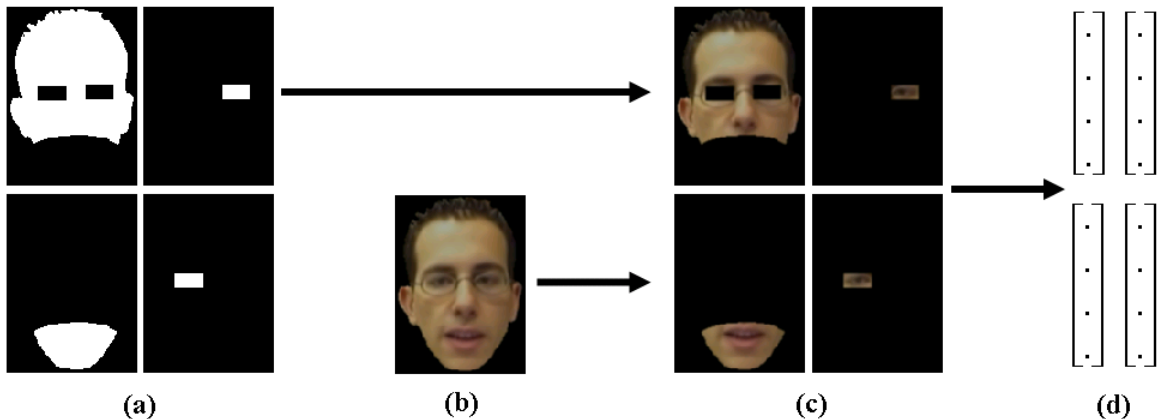


Figure 1: (a) Masks π^r . (b) Image \mathbf{I}_t . (c) Regions \mathbf{R}_t^r , obtained from the application of each mask π^r over image \mathbf{I}_t . (d) Vectors \mathbf{o}_t^r related to the defined regions.



Figure 2: (a) Three observed frames of a subject's face. (b) The same synthesized frames after learning the appearance model of this person.

the corresponding alignment coefficients \mathbf{a}_1 are set to $\mathbf{0}$; they represent the affine transformation used to fit the masks onto the face on each frame [5]. Using the tracking algorithm presented in [7] over the second image \mathbf{I}_2 , observations $\mathbf{o}_2^1, \dots, \mathbf{o}_2^4$ and alignment coefficients \mathbf{a}_2 are stored. At this point, each facial element r can be factorized as $[\mathbf{o}_1^r \ \mathbf{o}_2^r] = \mathbf{O}_2^r = \mathbf{U}_2^r \boldsymbol{\Sigma}_2^r (\mathbf{V}_2^r)^T + \bar{\mathbf{o}}_2^r = \mathbf{U}_2^r (\mathbf{C}_2^r)^T + \bar{\mathbf{o}}_2^r$, where the current mean observation is generated by $\bar{\mathbf{o}}_2^r = \frac{\mathbf{o}_1^r + \mathbf{o}_2^r}{2}$, the eigenvectors of \mathbf{O}_2^r (\mathbf{O}_2^r)^T are found in \mathbf{U}_2^r and the texture parametrization of the r -th facial element in images \mathbf{I}_1 and \mathbf{I}_2 is obtained in \mathbf{C}_2^r . Once this initialization is done, the *On-the-fly Training Algorithm* (Figure 3) can be executed. Besides, only those columns of \mathbf{U}_{t+1}^r and \mathbf{V}_{t+1}^r whose values of $\boldsymbol{\Sigma}_{t+1}^r$ exceed a threshold τ are considered, keeping only those eigenvectors with enough information. The value of τ decreases from 0, 5 to 0, 5 · 10⁻³ in the first images (1 seconds at 25 im/s) in order to allow better face localization when almost no information is known about its appearance [14]. Notice that alignment parameters \mathbf{a} can be used to extract gestural information in a multimodal input system [16].

On-the-fly Training Algorithm

In: $\mathbf{U}_2, \boldsymbol{\Sigma}_2, \mathbf{V}_2, \bar{\mathbf{o}}_2$, alignment coefficients \mathbf{a}_2 and set of four masks $\boldsymbol{\Pi}$

1. Set $k = 2$
2. Using $\mathbf{U}_k, \boldsymbol{\Sigma}_k, \mathbf{V}_k, \bar{\mathbf{o}}_k, \mathbf{a}_k$ and $\boldsymbol{\Pi}$, the images of the new update block are aligned, generating L observation vectors \mathbf{o}_{k+1}^r and alignment information \mathbf{a}_{k+1} for each image.
3. Obtain $\mathbf{U}_{k+1}, \boldsymbol{\Sigma}_{k+1}, \mathbf{V}_{k+1}$ and $\bar{\mathbf{o}}_{k+1}$ from $\mathbf{U}_k, \boldsymbol{\Sigma}_k, \mathbf{V}_k, \bar{\mathbf{o}}_k$, and \mathbf{o}_{k+1} (4)-(8).
4. Trim \mathbf{U}_{k+1} and \mathbf{V}_{k+1} according to $\boldsymbol{\Sigma}_{k+1}$.
5. Set $k = k + 1$ and go to Step 2 until there is no more new images.

Out: $\mathbf{U}_f, \boldsymbol{\Sigma}_f, \mathbf{V}_f, \bar{\mathbf{o}}_f$ and alignment coefficients \mathbf{a}_f for each image.

3.3 Cost analysis

In this section, the computational cost and memory requirements of the incremental computation of matrices $\mathbf{U}_f, \boldsymbol{\Sigma}_f$ and \mathbf{V}_f and vector $\bar{\mathbf{o}}_f$ of the previous *On-the-fly algorithm* is presented in table 2. As could be seen in the previous section, this incremental process consists of successive SVD and data mean updates, explained in section 2.3 to achieve a final factorization of the whole observation matrix \mathbf{O} :

$$\mathbf{O}_{p \times q} = (\mathbf{U}_f)_{p \times s} (\boldsymbol{\Sigma}_f)_{s \times s} \left((\mathbf{V}_f)_{q \times s} \right)^T + \bar{\mathbf{o}}_{p \times 1} \cdot \mathbf{1}_{1 \times q} \quad (16)$$

The value s consists in the number of eigenvectors kept in matrices \mathbf{U}_{k+1} (step 4 of the *On-the-fly algorithm*). Also, it must be noted that the update block size is specified by c . In this analysis, we presuppose that $p > q$,

Table 2: Resource order requirements of the proposed SVD and mean update algorithm over a matrix $\mathbf{O}_{p \times q}$ using update block size c and the s eigenvectors corresponding to the largest s singular values of $\mathbf{O}_{p \times q}$. To obtain more compact expressions, value $n = s + c$ has been used. Value k identifies the iteration number.

Id. Operation	Computational cost	Memory requirements
SVD $\left(\begin{bmatrix} \Sigma_k & \mathbf{U}_k^T \hat{\mathbf{I}} \\ \mathbf{0} & \mathbf{Q}_k^T \hat{\mathbf{I}} \end{bmatrix}_{n \times n} \right)$	$O(n^2)$	$O(n^2)$
$[\mathbf{U}_k \quad \mathbf{Q}_k]_{p \times n} \cdot (\mathbf{U}_d)_{n \times n}$	$O(pn^2)$	$O(pn)$
$(\mathbf{V}_d^T)_{n \times n} \cdot \begin{bmatrix} \mathbf{V}_k^T & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}_{n \times (kc+c)}$	$O((kc+c)n^2)$	$O((kc+c)n)$
Mean update	$O(s^2(p+kc+c))$	$O(s(p+kc+c))$
Total	$O\left(q\left(\frac{s^2}{c}+n\right)(n+p+q)\right)$	$O(n(p+q+s)+c^2)$

$q > s$ and $q > c$. Moreover, if additional considerations are taken into account for the values of c and s , particular cost functions can be described as follows:

- When c and s are of small order of magnitude (o.o.m.) compared to q , the lowest computational cost is obtained: $O(sq(p+q))$.
- If only c has small o.o.m., the computational cost becomes the highest one: $O(qs\frac{s}{c}(s+p+q))$.
- For small o.o.m of s only, the computational cost becomes $O(qc(c+p+q))$.
- When all c, s, p and q are of the same o.o.m., $O(q(s+c)(s+c+p+q))$.

The computational cost order of the batch process is $O(pq(p+q))$, which is higher than the first assumption and slightly higher than the two last ones. Note that the two last cases have also a similar cost.

Regarding to memory costs, the batch process has memory requirements of order $O(q^2+sp)$, while the proposed incremental approach has $O((c+s)(p+q+s)+c^2)$. As can be noted, for small values of c and s the presented approach achieves great memory reduction and do not increase its order in the other cases.

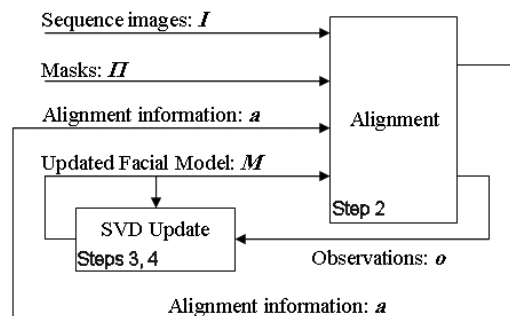


Figure 3: Block diagram of the *On-the-fly Training Algorithm*

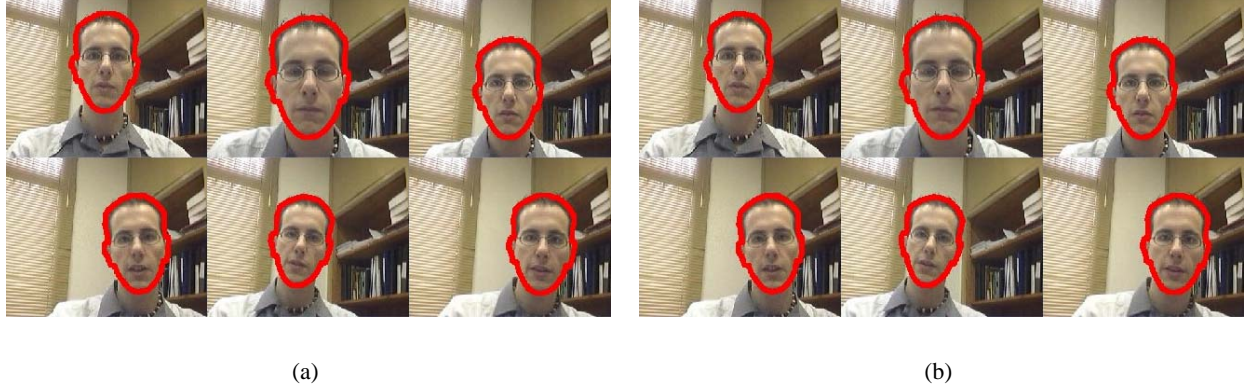


Figure 4: Output tracking results of the learning process for: (a) the *On-the-fly Training Algorithm* and (b) the non-causal algorithm.

4 Experimental Results

In this section, the performance of our incremental algorithm is shown. First, tracking results are specified in section 4.1, showing a comparison between the presented algorithm and its previous version. Next, precision results about incremental SVD with mean update are presented in section 4.2. Finally, in section 4.2.2 execution time is put in correspondence with the cost analysis obtained from section 3.3.

4.1 On-the-fly training algorithm

The *On-the-fly Training Algorithm* has been tested over a short sequence and a long one, both recorded at a frame rate of 25 im/s. The short sequence consists of 316 images and it has been used to compare the results obtained from our *On-the-fly Training Algorithm* and its previous non-causal version [7]. Achieving the same quality in the results (see Fig. 4), the presented algorithm has reduced the execution time about 66% with respect to [7] and has required about 7 Mbytes in front of the 200 Mbytes consumed by [7] (see the comparison in Fig. 5). Later, if we focus on the long sequence (10000 frames), its processing requirements were impossible to met with the non-causal algorithm [7] because its huge memory cost of 6000 Mbytes, although massive storage systems (e.g. hard drives) were used; the *On-the-fly Training Algorithm* reduced the memory requirements to 17 Mbytes with a processing time of a little more than 10 hours (using a 2GHz processor) (see Fig. 5).

4.2 Incremental SVD and mean computation

In this section, the goodness of the results given by the proposed incremental SVD and mean update algorithm (sect. 2.3) is analyzed and compared to the ideal performance offered by the batch solution.

4.2.1 Precision comparisons

Some experiments have been developed in order to test the analysis shown in the previous section (3.3). Two video sequences have been recorded and the face has been aligned in each one using our *On-the-fly training algorithm*. Starting from these aligned observations set stored columnwise in every \mathbf{O}^k , we have factorized it using both the batch SVD process and our incremental SVD with mean update algorithm (sect. 2.3), obtaining two approximations of the form:

$$\mathbf{O}_{p \times q}^k \approx \mathbf{U}_{p \times s}^k \mathbf{\Sigma}_{s \times s}^k \left(\mathbf{V}_{q \times s}^k \right)^T + \bar{\mathbf{o}}_{p \times 1}^k \cdot \mathbf{1}_{1 \times q} \quad (17)$$

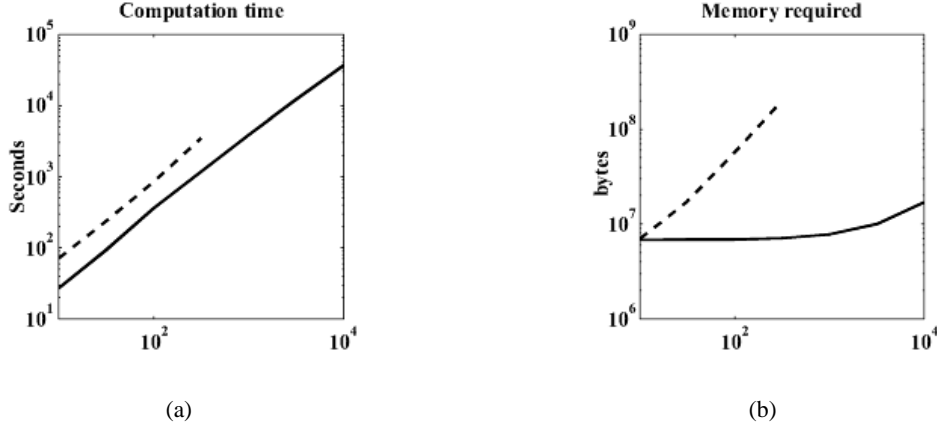


Figure 5: Solid line represents the *On-the-fly Training Algorithm* performance while the dashed one belongs to the non-causal algorithm presented in [7]. (a) Computation time in seconds. (b) Memory used in bytes.

$$\mathbf{O}_{p \times q}^k \approx \hat{\mathbf{U}}_{p \times s}^k \hat{\Sigma}_{s \times s}^k \left(\hat{\mathbf{V}}_{q \times s}^k \right)^T + \hat{\mathbf{o}}_{p \times 1}^k \cdot \mathbf{1}_{1 \times q} \quad (18)$$

where matrices \mathbf{U}^k , Σ^k and \mathbf{V}^k are the trimmed version of the thin-SVD of $\hat{\mathbf{O}}^k = \mathbf{O}^k - \bar{\mathbf{o}}^k \cdot \mathbf{1}$ and $\bar{\mathbf{o}}^k$ is the mean column of \mathbf{O}^k ; matrices $\hat{\mathbf{U}}^k$, $\hat{\Sigma}^k$ and $\hat{\mathbf{V}}^k$ and vector $\hat{\mathbf{o}}^k$ are the corresponding ones when obtained with the incremental approach presented in section 2.3. This incremental process has been executed with different sizes of update block c and different threshold τ (sect. 3.2); the higher the threshold, the lesser eigenvalues kept in the model (with a non-linear case specific relation $s = f(c, \tau, k)$). Next, we define:

$$e_b(c, \tau) = \sum_{\forall k} \left\| \mathbf{M}_{p \times q}^k - \mathbf{U}_{p \times s}^k \Sigma_{s \times s}^k \left(\mathbf{V}_{q \times s}^k \right)^T - \bar{\mathbf{o}}_{p \times 1}^k \cdot \mathbf{1}_{1 \times q} \right\|_2 \quad (19)$$

$$e_i(c, \tau) = \sum_{\forall k} \left\| \mathbf{M}_{p \times q}^k - \hat{\mathbf{U}}_{p \times s}^k \hat{\Sigma}_{s \times s}^k \left(\hat{\mathbf{V}}_{q \times s}^k \right)^T - \hat{\mathbf{o}}_{p \times 1}^k \cdot \mathbf{1}_{1 \times q} \right\|_2 \quad (20)$$

Function e_b is shown in fig. 6(a) and e_i is represented in fig. 6(b). Following the reduction and compression of matrices theorem found in [15], it can be assured that $e_b(c, \tau) \leq e_i(c, \tau)$ for any c and τ . Figure 6(c) represents the relative error as a function of c and τ . This relative error is measured as $\frac{e_i(c, \tau) - e_b(c, \tau)}{e_b(c, \tau)}$ and, as can be observed, all three figures achieve it lowest value when both c and τ have low values (1-5 and 0.001, respectively).

4.2.2 Execution time

We have measured the execution time of both the batch and our incremental computation process done in section 4.2.1. The execution time of our incremental SVD and mean update algorithm is depicted in fig. 7 as a function of update block size and treshold (sect. 3.2) and has been obtained as the mean execution time related to the observation matrices \mathbf{O}^k .

It can be noted that the analysis made in sect. 3.3 is reflected in fig. 7. It must be noted that the fastest results (about a third of the computation time belonging to the batch approach) can be achieved for small update block sizes and large threshold, which translates in taking into account few (1-2) eigenvectors. By the other hand, the heaviest computational load corresponds to the assumption of small block size (1-5) and low threshold (0.001), which translates to a larger number of eigenvectors (30) and further overcomes the computation time of the batch process. Finally, it can also be seen that as the block size grows, the computational cost becomes more independent with respect to the threshold (or number of eigenvector kept).

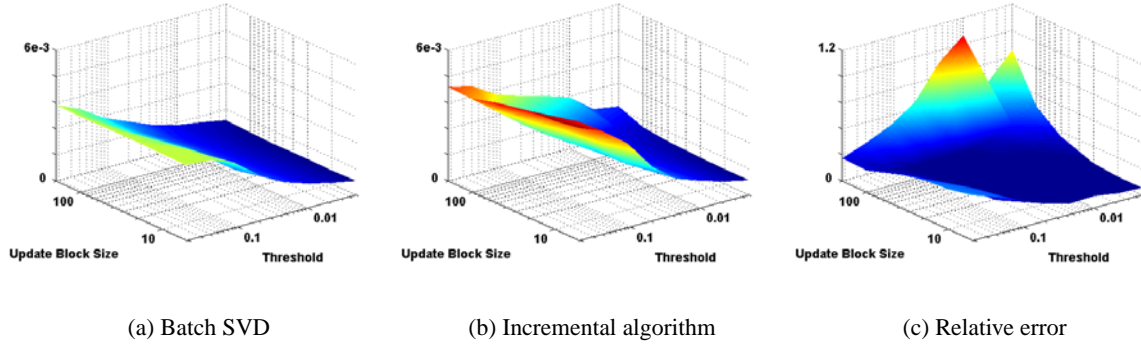


Figure 6: (b) Error between the original data matrix \mathbf{O} and the factorization obtained by our incremental SVD and mean update algorithm.

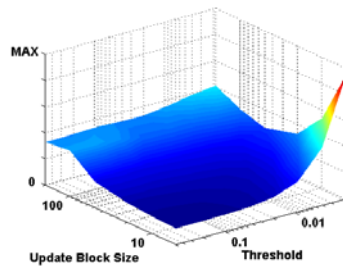


Figure 7: Execution time of the algorithm with different number of eigenvectors and update block size

4.2.3 Conclusions

It can be concluded that the best alternative consists in using a small block size (*i.e.* 1-10) with a relatively small threshold (*i.e.* 0.01, obtaining about 10 eigenvectors); it achieves a relative error of less than 10^{-3} with half the computation time of the corresponding batch process. Moreover, when both update block size and threshold are small enough ($\tau = 0.001$, obtaining more than 30 eigenvectors, and $c = 1$), the incremental SVD and mean update algorithm achieves the best performance but with the heaviest computational load. By the other hand, the fastest option, achieved with small update block size and high threshold ($\tau = 0.1$, $c = 1$), offers a poor precision compared to the previous cases. Finally, if we increase the update block size ($c > 10$), both computational and precision results also get worse.

5 Concluding Remarks

In this paper, a new method for extracting the mean of an existing SVD is presented, without increasing either the cost order of memory or time. This fact has allowed us to offer an incremental computation of SVD preserving a zero data mean, which has been analyzed and compared with the batch approach. The precision offered by our method is high enough to allow photorealistic reconstructions of observed face images using half the computation time of the non-incremental processes. Fields that can benefit from it can be, *e.g.*: classification problems, where the mean information is used to center the data; incremental computation of covariation matrices, which need to be centered around its mean; causal construction of eigenspaces, where the principal components of the data are included, as well as the mean information. With respect to the latter, the *On-the-fly Algorithm* is presented in this work. Given an image sequence and a set of masks, this algorithm is capable of generating a separate eigenspace for each facial element (learning all their appearance variations due to changes

in expression and visual utterances) and effectively tracking and aligning them. Furthermore, longer sequences than previous methods [5], [7] can be processed with the same visual accuracy when no illumination changes appear. Finally, we plan to add more robustness to this algorithm using methods like [5] and more work will be done in order to achieve real time performance, so specific appearance models can be obtained as a person is being recorded.

References

- [1] E. André: The generation of multimedia presentations. *A Handbook of Natural Language Processing: techniques and applications for the processing of language as text*, R. Dale, H. Misl and H. somers, Eds., Marcel Dekker Inc., 2000, 305-327
- [2] S. Robbe-Reiter, N. Carbonell and P. Dauchy: Expression constraints in multimodal human-computer interaction. *Intelligent User Interfaces*, 2000, 225-228
- [3] C. Stephanidis: Towards universal acces in the disappearing computer environment. *UPGRADE*, vol. IV, n. 1, 2003, pp. 53-59
- [4] E. André: Natural language in multimedia/multimodal systems. *Handbook of Computational Linguistics*, R. Miktov, Ed., Oxford Univ. Press, 2003, 650-669
- [5] F. de la Torre and M. Black: Robust parameterized component analysis: Theory and applications to 2d facial modeling. *ECCV*, 2002, 654-669
- [6] T. Ezzat, G. Geiger and T. Poggio: Trainable videorealistic speech animation. *ACM SIGGRAPH*, San Antonio, Texas, July 2002, 225-228
- [7] J. Melenchón, F. de la Torre, I. Iriondo, F. Alías, E. Martínez and Ll. Vicent: Text to visual synthesis with appearance models, *ICIP*, 2003, vol I, 237-240
- [8] D. Cosker, D. Marshall, P. Rosin and Y. Hicks: Video realistic talking heads using hierarchical non-linear speech-appearance models. *Mirage*, France, 2003
- [9] B.J. Theobald , J.A. Bangham, I. Matthews and G.C. Cawley: Near-videorealistic synthetic talking faces: Implementation and evaluation (submitted on invitation), *Speech Communication Journal*, 2004.
- [10] M. Gu and S.C. Eisenstat: A Stable and fast algorithm for updating the singular value decomposition. *Tech. Rep. YALEU/DCS/RR-966*, New Haven, 1993
- [11] S. Chandrasekaran, B. Manjunath, Y. Wang, J. Winkeler and H. Zhang: An eigenspace update algorithm for image analysis. *GMIP*, vol. 59, no. 5, 1997, 321-332
- [12] M. Brand: Incremental singular value decomposition of uncertain data with missing values. *ECCV*, 2002, I: 707 ff.
- [13] P.M. Hall, D.R. Marshall and R. Martin: Adding and subtracting eigenspaces with eigenvalue decomposition and singular value decomposition. *ICV*, vol. 20, no. 13-14, december 2002, 1009-1016
- [14] D. Scočaj: Robust Subspace Approaches to Visual Learning and Recognition. Ph.D. dissertation, University of Ljubljana, Faculty of computer and information science, Ljubljana, 2003.
- [15] M. Kirby: *Geometric Data Analysis: An Empirical Approach to Dimensionality Reduction and the Study of Patterns*. John Wiley & Sons Inc., 2001

- [16] F. Keates and P. Robinson: Gestures and multimodal input. *Behaviour and Information Technology*, Taylor and Francis Ltd., 1999, 18(1), 35-42