

DICTIONARY LEARNING FOR A SPARSE APPEARANCE MODEL IN VISUAL TRACKING

Sylvain Rousseau¹, Pierre Chainais^{1,2}

¹ École Centrale de Lille
CRIStAL, UMR CNRS 9189, France
² INRIA Lille-Nord Europe, Sequel
Villeneuve d’Ascq, France

sylvain.rousseau@ec-lille.fr

pierre.chainais@ec-lille.fr

Christelle Garnier

Institut Mines-Telecom / Telecom Lille
CRIStAL, UMR CNRS 9189, France

christelle.garnier@telecom-lille.fr

ABSTRACT

This paper presents a novel approach to visual object tracking based on particle filtering. The tracked object is modelled by a sparse representation provided by dictionary learning. Such an approach permits to describe the target by a model of reduced dimension. The likelihood of a candidate region is built on a similarity measure between the sparse representations of a set of patches (at known positions) in the dictionary learnt from the reference template. Experimental validation is performed on various video sequences and shows the robustness of the proposed approach.

Index Terms— dictionary learning, sparse coding, particle filtering, object tracking

1. INTRODUCTION

Tracking objects in a video sequence is an essential task in computer vision. Among the many tracking methods (see [1]), particle filters are currently widely used. The performance strongly depends on the appearance model used to characterise the target object and on the likelihood function which measures the similarity between the candidate region and the reference object. Recently new models based on sparse representation, initially introduced in [2], have been successfully applied in visual tracking [3]. A global sparse appearance model is proposed in [4]. The dictionary is composed of two parts: target templates used to represent the entire object and trivial templates introduced to account for noisy pixels. A candidate is then sparsely represented with a few atoms of the complete dictionary. The likelihood is related to the reconstruction error in the target dictionary. The discriminative ability of the method can be improved by adding background templates in the dictionary [5]. To better handle partial occlusions, a local sparse appearance model based on a dictionary of patches sampled from target

templates [6] or learnt from these patches [7] is more appropriate. The principle is to extract patches from a candidate and to compute their sparse encoding in the dictionary. The representation coefficients or sparse codes are then exploited as appearance features. The eventual step, referred to as the pooling step, is to extract meaningful coefficients from these appearance features. The likelihood is based on the similarity between the candidate and the object features vectors.

Several pooling operators have been considered. The average pooling [8] averages the sparse codes to build up the final feature vector. The max pooling method [7] keeps the maximum absolute value of the sparse code coefficients. Both methods are imprecise because they lose the location of patches and thus the spatial structure of the object. To overcome this loss, an alignment pooling method is developed in [9]. It preserves the spatial information, but it imposes constraints on the dictionary construction. It is composed of patches extracted from target templates according to a fixed spatial structure that can not change over time.

In this paper we propose to combine the principle of alignment pooling with dictionary learning for sparse representation of the target object patches. Then tracking is performed in a space of reduced dimension by comparing sparse codes of patches in the candidate region with sparse codes of patches extracted from the target template at known locations. The paper is organised as follows. Section 2 recalls related previous works. Section 3 presents the proposed sparse coding model. In Section 4, the tracking algorithm based on particle filtering is described. Section 5 presents performance results obtained on public datasets and Section 6 concludes.

2. RELATED WORK

We first recall the alignment pooling approach [9]. The target is described by several templates. For simplicity, we consider only one target template T , typically a rectangle comprising M overlapping square patches of size c^2 . A set of $N \leq M$ patches is extracted from the template at predefined

Thanks to the BNPSI ANR project no ANR-13-BS-03-0006-01, the AR-CIR REPAR project and to the Fondation École Centrale Lille for funding.

locations according to a spatial layout. The matrix $\mathbf{L}_p = [p_1, \dots, p_N] \in \mathbb{R}^{c^2 \times N}$, formed by the concatenation of the vectorised patches $(p_i)_{i=1}^N$ is directly used as the dictionary. Then a candidate region C is modelled as follows. The candidate is first resized to match the target template's size. A set of N patches $(q_i)_{i=1}^N$ is sampled at the *same locations* as in the target. The sparse codes \mathbf{u}_i of each of these patches in the dictionary \mathbf{L}_p of size N are obtained by solving the following ℓ_1 -constrained minimisation problem

$$\arg \min_{\mathbf{u}_i} \|\mathbf{q}_i - \mathbf{L}_p \mathbf{u}_i\|_2^2 \quad \text{s.t.} \quad \|\mathbf{u}_i\|_1 \leq \rho, \quad (1)$$

where ρ controls the sparsity of the code $\mathbf{u}_i \in \mathbb{R}^N$. Each patch \mathbf{q}_i is then approximated by a sparse combination of the patches $(p_i)_{i=1}^N$ weighted by \mathbf{u}_i . The alignment pooling of \mathbf{u}_i consists in selecting the contribution of patch p_i at the same location as patch \mathbf{q}_i which is the i -th coefficient u_{ii} of \mathbf{u}_i . Consequently, if the \mathbf{u}_i 's are the columns of the square matrix $\mathbf{U} \in \mathbb{R}^{N \times N}$, the pooled coefficients are located on the diagonal. The log-likelihood of candidate C with respect to template T is defined as the trace of matrix \mathbf{U}

$$\mathcal{L}(C, T) = \text{Tr} \mathbf{U} = \sum_{i=1}^N u_{ii}. \quad (2)$$

For comparison, the max and average pooling operators respectively compute the maximum and average of each \mathbf{u}_i , before summing the results. The alignment pooling is more accurate because it takes into account the locations of extracted patches. However the major drawback is the lack of flexibility. To increase the number of patches or change their locations for better performance, the dictionary has to be completely redefined. The size of the dictionary is necessarily the number N of extracted patches so that if one wants to take all available patches into account, the size of the dictionary becomes large since $N = M$. Moreover, as shown on Fig. 2, that method may suffer from a lack of robustness.

3. PROPOSED APPROACH

To overcome these limitations, we propose to learn a dictionary $\mathbf{D}_p \in \mathbb{R}^{c^2 \times K}$ of fixed size $K \ll M$ from a set of patches $\mathbf{L}_p = [p_1, \dots, p_N] \in \mathbb{R}^{c^2 \times N}$ extracted from template T to obtain an adapted representation of the data. This is carried out by solving the minimisation problem

$$\arg \min_{\mathbf{D}_p, \mathbf{Z}} \|\mathbf{L}_p - \mathbf{D}_p \mathbf{Z}\|_F^2 \quad \text{s.t.} \quad \|\mathbf{z}_i\|_1 \leq \rho, \quad \forall i, \quad (3)$$

where \mathbf{z}_i is the i -th column of the matrix $\mathbf{Z} \in \mathbb{R}^{K \times N}$ and $\|\cdot\|_F$ denotes the Frobenius norm. Dictionary learning is performed on centered patches while average values are processed apart. Consequently, \mathbf{z}_i is the sparse code of the patch p_i with respect to the dictionary \mathbf{D}_p . The dimension reduction performed by dictionary learning imposes no limit

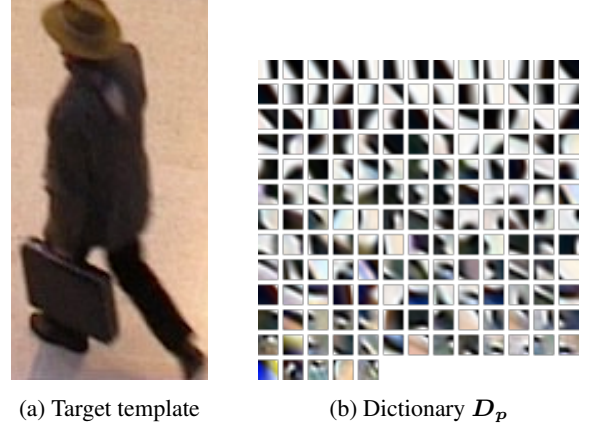


Fig. 1: Atoms of a learnt dictionary by order of contributions in the image.

on the number N . We can afford to use all the overlapping patches ($N = M$) of size c^2 within the template to learn the dictionary of size $K \ll M$. Figure 1 illustrates the atoms of the dictionary \mathbf{D}_p of size $K = 161$ learnt over all the $N = 16109$ patches of size 8×8 extracted from the target template image in Fig. 1a. To compare, the alignment pooling has to solve the minimisation problem (1) with a dictionary of size 16109 whereas the proposed method uses a dictionary \mathbf{D}_p that is 100 times smaller. To describe a candidate region, n patches $(q_i)_{i \in J}$ are extracted at locations defined by some set $J \subset \{1, \dots, N\}$ and $\#J = n$. We compute their sparse codes $\mathbf{V} = [\mathbf{v}_i]_{i \in J} \in \mathbb{R}^{K \times n}$ in \mathbf{D}_p by solving

$$\arg \min_{\mathbf{v}_i} \|\mathbf{q}_i - \mathbf{D}_p \mathbf{v}_i\|_2^2 \quad \text{s.t.} \quad \|\mathbf{v}_i\|_1 \leq \rho. \quad (4)$$

Compared to [9] where $n = N$, the choice of the number of patches n and their locations is completely free. To build the log-likelihood of a candidate region C , we compare the sparse representations using \mathbf{D}_p of these n patches from C and T at the same locations defined by J . This comparison is quantified by matched filtering which provides optimal SNR for pattern detection. Thus, the log-likelihood of a candidate set of n sparse codes \mathbf{V} with respect to the template sparse codes $\mathbf{Z}_J = [\mathbf{z}_i]_{i \in J}$ is defined as a cross-correlation:

$$\frac{1}{\|\mathbf{Z}_J\|_F} \text{Tr} \mathbf{Z}_J^T \mathbf{V}. \quad (5)$$

This log-likelihood optimally detects the presence of the template patches represented by \mathbf{Z}_J at known positions. This approach can be interpreted as an alignment pooling strategy since the *trace* operation pools the individual alignment scores $\mathbf{z}_i^T \mathbf{v}_i$ of the set of n patches (which all ‘vote’ for the same target position). Note that the similarity proposed in [9] can be seen as a special case of (5) where $\mathbf{D}_p = \mathbf{L}_p$ (no dictionary learning) and $\mathbf{Z}_J = \mathbf{I}$ (trivial coding). Last remains to take into account the average values of patches denoted by

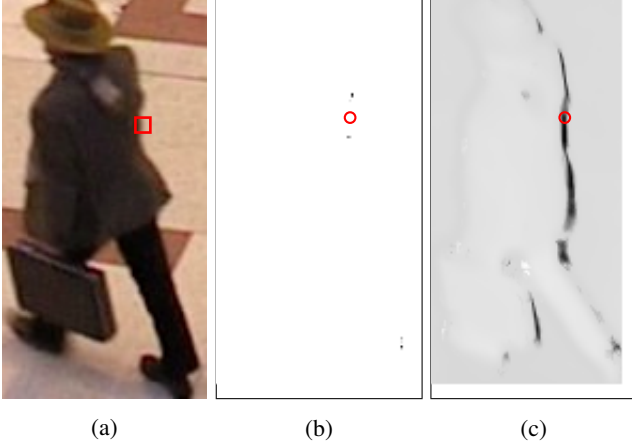


Fig. 2: (a) A patch q_{i_0} taken from a resized candidate. (b) Reshaped sparse code u_{i_0} . (c) Cross-correlation $Z_J^T v_{i_0}$

\bar{p}_i and \bar{q}_i : large differences of average gray levels $|\bar{p}_i - \bar{q}_i|$ should be penalised. We propose to add a term $-\lambda|\bar{p}_i - \bar{q}_i|^2$ to (5) to finally get the log-likelihood

$$\mathcal{L}_J(C, T) = \frac{1}{\|Z_J\|_F} \text{Tr} Z_J^T V - \lambda \sum_{i \in J} |\bar{p}_i - \bar{q}_i|^2, \quad (6)$$

where $\lambda > 0$ is some weighting parameter (default is 1).

Figure 2 compares the proposed approach to the alignment pooling method. A patch q_{i_0} is selected in a resized candidate region of Fig. 2a. Figures 2b and 2c show the similarity maps between q_{i_0} and all the patches $(p_i)_{i \in J}$ with $J = \{1, \dots, M\}$ extracted from T . Without dictionary learning, see Fig. 2b, this similarity map is directly given by the sparse code u_{i_0} of the patch q_{i_0} in L_p . With the proposed approach at Fig. 2c, it is given by $Z_J^T v_{i_0}$ where v_{i_0} is the sparse code of q_{i_0} in D_p . The alignment pooling consists in selecting the coefficient of the sparse code corresponding to the contribution of the patch p_{i_0} at the same location as q_{i_0} (red circles in Fig. 2b and 2c). With a standard alignment pooling, little matching is found between the patches p_{i_0} and q_{i_0} . In the set of the template patches used as the dictionary, a lot of atoms (patches) are similar and the atoms selected by the sparse decomposition do not necessarily include the aligned patch p_{i_0} . This effect can be reduced by considering templates from different frames. Our approach gives better matching results because the learnt dictionary is of reduced dimension and the template patches $(p_i)_{i \in J}$ are described by several atoms of D_p . A match is found every time $p_i, i \in J$, shares atoms with q_{i_0} . Figure 2c shows several matches in particular with the aligned patch.

4. TRACKING ALGORITHM

Object tracking is carried out by particle filtering [10]. Given a sequence of observations $\mathbf{y}_{1:k} = (\mathbf{y}_1, \dots, \mathbf{y}_k)$ up to time k ,

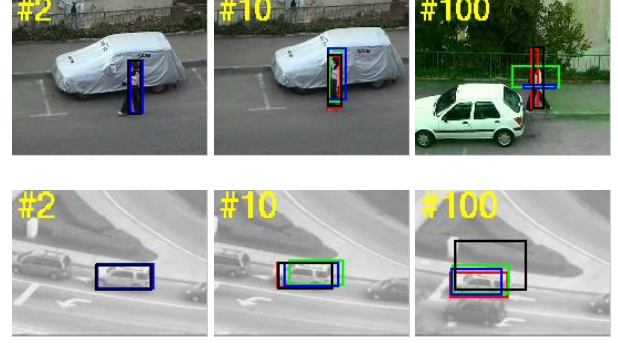


Fig. 3: Qualitative tracking results on the sequences Woman and SUV with different pooling strategies: proposed in red, maximum pooling in green, average pooling in blue and alignment pooling in black.

the aim is to estimate the object state \mathbf{x}_k from the posterior density $p(\mathbf{x}_k | \mathbf{y}_{1:k})$ which is recursively approximated with a set of N_p weighted particles $\{\mathbf{x}_k^{(m)}, w_k^{(m)}\}_{i=1}^{N_p}$: $p(\mathbf{x}_k | \mathbf{y}_{1:k}) \approx \sum_{i=1}^{N_p} w_k^{(m)} \cdot \delta(\mathbf{x}_k - \mathbf{x}_k^{(m)})$. The target object is represented by a bounding window. The state vector is defined as $\mathbf{x}_k = [c_k, d_k]^T$ where $c_k = [c_k^x, c_k^y]$ is the position of the top left corner and $d_k = [d_k^x, d_k^y]$ is the size of the window. To obtain the set of particles at time k from the previous particles $\{\mathbf{x}_{k-1}^{(m)}, w_{k-1}^{(m)}\}_{i=1}^{N_p}$, the particle filter algorithm includes two major steps: prediction and update. During prediction, the samples $\mathbf{x}_k^{(m)}$ are propagated according to the dynamic model using the prior density $p(\mathbf{x}_k^{(m)} | \mathbf{x}_{k-1}^{(m)})$. As often, we assume that \mathbf{x}_k evolves as a Gaussian random walk: $\mathbf{x}_k | \mathbf{x}_{k-1} \sim \mathcal{N}(\mathbf{x}_{k-1}, \Sigma)$ where Σ is a diagonal covariance matrix which defines the uncertainty region around the previous state. During the second stage, the weights $w_k^{(m)}$ are updated according to the recursive expression: $w_k^{(m)} = w_{k-1}^{(m)} \cdot p(\mathbf{y}_k | \mathbf{x}_k^{(m)})$ where the observation likelihood $p(\mathbf{y}_k | \mathbf{x}_k)$ measures the matching between the observation and the state. We define the likelihood directly from expression (6) as

$$p(\mathbf{y}_k | \mathbf{x}_k^{(m)}) \propto \exp \{ \mu \cdot \mathcal{L}_J(C, T) \}, \quad (7)$$

where C is the candidate region identified by $\mathbf{x}_k^{(m)}$ and μ is a tuning parameter. Finally, after weight normalisation and if necessary, particle resampling, the object state is obtained by MMSE (Minimum Mean Square Error) estimation

$$\widehat{\mathbf{x}}_k = \mathbb{E} [\mathbf{x}_k | \mathbf{y}_{1:k}] = \sum_{m=1}^{N_p} w_k^{(m)} \mathbf{x}_k^{(m)}. \quad (8)$$

5. EXPERIMENTAL VALIDATION

In this section, we validate our approach by comparing it to other pooling strategies. The optimisation problems (1), (3)

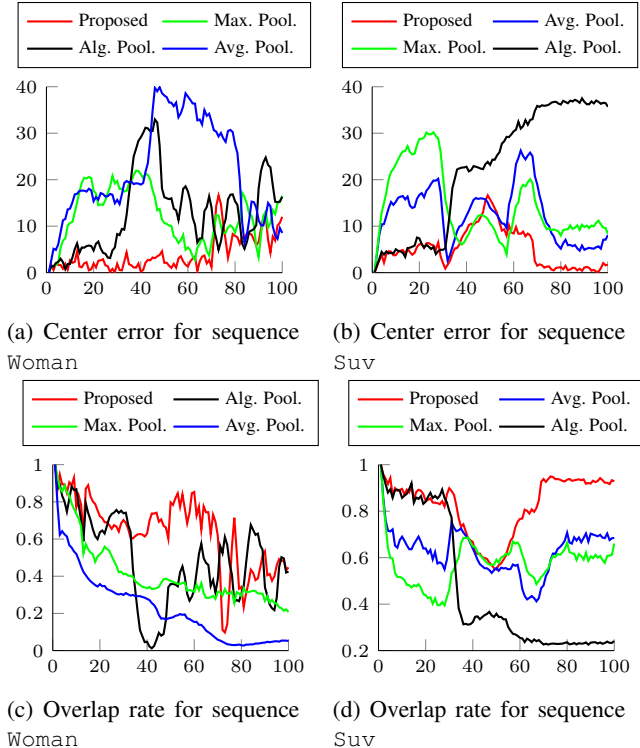


Fig. 4: Quantitative tracking results on the sequences *Woman* and *Suv*

and (4) are solved using the Spams Matlab toolbox¹ [11] with $\rho = 0.25$. J is set to a fixed grid of step 8. For particle filtering, $N_p = 600$ particles, $\mu = 4.6$ and $\Sigma = \text{diag}(20, 20, 4, 4)$. Concerning the proposed method, the size of the dictionary is set to $K = 64$.

Four pooling strategies are considered: alignment pooling (without dictionary learning) and maximum, average or alignment pooling (proposed) with dictionary learning. They are compared in the first 100 frames of the sequences *Suv* and *Woman*. For a fair comparison, all the trackers adopt the same dynamical model, only the observation likelihood differs. Figure 3 shows the qualitative tracking results obtained on a few images of the sequences. These sequences present challenging scenarios with partial occlusion and illumination changes. Figure 4 shows the quantitative tracking results. The first row depicts the center location error (in pixels) between the ground truth bounding box and the estimated bounding box. The second row depicts the overlap rate [12] which is the ratio between the area of the intersection of both bounding boxes and the area of the union.

The maximum and average pooling give the worst performance in the first frames of both videos, because they use a similarity between patches regardless of their locations. However this lack of precision is compensated by a good robust-

	$R(\alpha = 0.5)$			AUC		
	Max	Alg	Prop	Max	Alg	Prop
Bolt	0.05	0.09	0.03	0.05	0.14	0.04
CarScale	0.14	1.00	0.98	0.23	0.74	0.74
David3	0.04	0.24	0.29	0.05	0.21	0.22
FaceOcc1	0.22	0.18	1.00	0.29	0.27	0.82
FleetFace	0.29	0.55	0.62	0.40	0.49	0.48
Freeman4	0.03	0.01	0.03	0.02	0.08	0.12
Singer1	0.34	0.59	0.61	0.29	0.57	0.65

Table 1: Success rate $R(\alpha = 0.5)$ and Area Under the Curve (AUC) $\int R(\alpha)d\alpha$ for the proposed method, max pooling and alignment pooling.

ness, and the performances are maintained or even can improve along the sequences. On the opposite, the alignment pooling is more accurate at the beginning of the sequences, because it takes into account the spatial information and a similarity is found between aligned patches. But the tracking results degrade along the videos. As shown in Figure 2, it can miss the matching between patches that are similar. For example on the *Suv* sequence, the target partially exits the image around frame 30. When the target comes back, the maximum and average pooling can find it whereas the alignment pooling fails. Along both sequences, the best tracking performance is achieved with our approach. The proposed cross-correlation based similarity is both accurate and robust. As the maximum and average pooling, our tracker can retrieve the target after a partial exit and as the alignment pooling, it preserves the spatial information.

Table 1 gathers results obtained on a set of publicly available video sequences [13]. It compares success rates $R(\alpha)$ that are the proportions of frames for which the overlap rate is $\geq \alpha$. To this aim, we compare the values of $R(\alpha = 0.5)$ and, in the spirit of the usual AUC, the Area Under the Curve $\{\alpha \in (0, 1), R(\alpha)\}$ which is a more general indicator. The proposed method either equals or outperforms other methods.

6. CONCLUSIONS

In this paper we propose a sparse representation of patches in a learnt dictionary to model the appearance of an object. Then tracking is performed by a matched filter on the sparse codes of the patches extracted from the candidate. This method generalises the alignment pooling method and improves on other pooling strategies with learnt dictionaries such that the max or average pooling. Numerical results show that it compares favourably with previous works.

¹<http://spams-devel.gforge.inria.fr/>

7. REFERENCES

- [1] H. Yang, L. Shao, F. Zheng, L. Wang, and Z. Song, "Recent advances and trends in visual tracking: A review," *Neurocomputing*, vol. 74, no. 18, pp. 3823–3831, 2011.
- [2] B. A. Olshausen and D. J. Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images," *NATURE*, vol. 381, p. 13, 1996.
- [3] S. Zhang, H. Yao, X. Sun, and X. Lu, "Sparse coding based visual tracking: Review and experimental comparison," *Pattern Recognition*, vol. 46, no. 7, pp. 1772–1788, 2013.
- [4] X. Mei and H. Ling, "Robust visual tracking using ℓ_1 minimization," in *IEEE International Conference on Computer Vision (ICCV)*, pp. 1436–1443, IEEE, 2009.
- [5] W. Zhong, H. Lu, and M.-H. Yang, "Robust object tracking via sparsity-based collaborative model," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1838–1845, IEEE, 2012.
- [6] B. Liu, J. Huang, L. Yang, and C. Kulikowsk, "Robust tracking using local sparse appearance model and k-selection," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1313–1320, IEEE, 2011.
- [7] Q. Wang, F. Chen, J. Yang, W. Xu, and M.-H. Yang, "Transferring visual prior for online object tracking," *IEEE Transactions on Image Processing*, vol. 21, no. 7, pp. 3296–3305, 2012.
- [8] S. Zhang, H. Yao, and S. Liu, "Robust visual tracking using feature-based visual attention," in *IEEE International Conference on Acoustics Speech and Signal Processing (ICASSP)*, pp. 1150–1153, IEEE, 2010.
- [9] X. Jia, H. Lu, and M.-H. Yang, "Visual tracking via adaptive structural local sparse appearance model," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1822–1829, IEEE, 2012.
- [10] M. Isard and A. Blake, "Condensation—conditional density propagation for visual tracking," *International journal of computer vision*, vol. 29, no. 1, pp. 5–28, 1998.
- [11] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, "Online learning for matrix factorization and sparse coding," *The Journal of Machine Learning Research*, vol. 11, pp. 19–60, 2010.
- [12] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [13] Y. Wu, J. Lim, and M.-H. Yang, "Online object tracking: A benchmark," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2411–2418, IEEE, 2013.