

Learning Rotation for Kernel Correlation Filter

Abdullah Hamdi, Bernard Ghanem

King Abdullah University of Science and Technology (KAUST), Thuwal, Saudi Arabia

{abdullah.hamdi, Bernard.Ghanem} @kaust.edu.sa

August 15, 2017

Abstract

Kernel Correlation Filters have shown a very promising scheme for visual tracking in terms of speed and accuracy on several benchmarks. However it suffers from problems that affect its performance like occlusion, rotation and scale change. This paper tries to tackle the problem of rotation by reformulating the optimization problem for learning the correlation filter. This modification (RKCF) includes learning rotation filter that utilizes circulant structure of HOG feature to guess rotation from one frame to another and enhance the detection of KCF. Hence it gains boost in overall accuracy in many of OBT50 detest videos with minimal additional computation

1 Introduction

Visual object tracking is a very important task in computer vision in which an object of interest would be identified and located in the first frame of a video. The goal is to follow the object movement and scale in subsequent frames by applying the tracking algorithm, usually faster and more efficient than a general detection scheme. Because of the wide range of applications that include visual tracking (e.g. robotics and surveillance), tracking had the attention of the computer vision community for several years

Correlation filters (CF) have used in tracking for several years in visual tracking due to their speed and efficient computations. Adding Kernels to these trackers produced state of the art Kernel Correlation Filter (KCF) that topped tracking benchmarks for several years. KCF utilizes the

circulant structure of the data matrix of all possible shifts to achieve less computation and utilizes the Kernel substitution trick [3]

Several versions and modifications of KCF came to tackle the problems it possessed like occlusion, boundary effect, scale, and rotation [4] [1][6][2].

Rotation being one of these problems that KCF suffers is itself an interesting problem (rotation detection) with wide range of applications like texture classification [5]. By the way KCF filter is constructed it assumes the object didn't rotate or change shape, this assumption cause the response of the filter to deteriorate and the detection would (as a result) drift away from the target position and cause drop of the performance.

We propose here to reformulate the optimization objective to include a second filter that will learn a rotation descriptor for the target and utilize this information in the detection phase of KCF family tracker. We show that this extra information (the rotation from one frame to another) will enhance the performance of CF trackers with many of its variations, with potential of applying in trackers that utilize deep features and Siamese network as discussed in [6]. This rotation filter uses the circularity of the HOG feature that enable it to utilize the same computational efficiency of the KCF tracker (as we will show in section 3) giving boost to the base performance with almost no additional computation time.

2 Related work

2.1 Kernel Correlation Filter

State of the art tracking technique that utilizes the cyclic nature of the shifted patches and Kernel trick to implement a very fast tracking algorithm based on correlation filters. They are filters that try to guesstimate the new position of the target by learning filters on each patch with expected response to be Gaussian with maximum at the center if the object didn't move .Figure 1 shows a typical response/target of regression of the filter

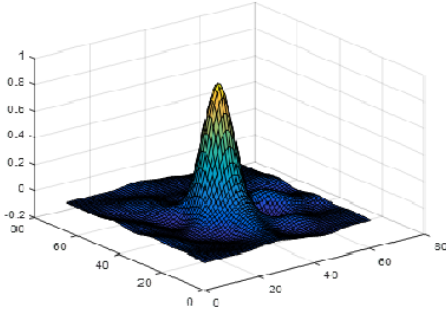


Figure 1: typical response of applying the learned KCF filter on an image patch

If the object moved by little from frame to another , the maximum response will be shifted and the translation of the maximum will be used to translate the patch in the image , and so forth in the following frames. [3]. To achieve this ,a filter w that minimize the energy of error between the response of the filter on patch image and a typical response of stationary patch according to the following equation

$$\min_w \| \mathbf{X}w - \mathbf{y} \|_2^2 + \lambda \| w \|_2^2 \quad (1)$$

where \mathbf{y} is Gaussian response of the filter if the patch didn't move and \mathbf{X} is the data matrix of the image patch that contains the target being tracked , λ is the regularization variable.

the closed form solution of the optimization 1 given by :

$$w = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (2)$$

It can be seen that the matrix \mathbf{X} is a circulant matrix of all possible shifts of the vector \mathbf{x} which is the vectorized image patch surrounding the target being tracked. Since the matrix \mathbf{X} is circulant it can diagonalized by the DFT matrix as follows :

$$\mathbf{X} = \mathbf{F} \text{diag}(\hat{\mathbf{x}}) \mathbf{F}^H \quad (3)$$

where $\hat{\mathbf{x}}$ is DFT of vector \mathbf{x} . Using 3 we can see that

$$\mathbf{X}^H \mathbf{X} = \mathbf{F} \text{diag}(\hat{\mathbf{x}}^* \odot \hat{\mathbf{x}}) \mathbf{F}^H \quad (4)$$

substituting 4 into 2 we get the following closed form solution to the learned filter \hat{w} in the Fourier domain in each frame.

$$\hat{w} = \frac{\hat{\mathbf{x}} \odot \hat{\mathbf{y}}}{\hat{\mathbf{x}} \odot \hat{\mathbf{x}} + \lambda} \quad (5)$$

The hat indicate DFT of the term. The solution in dual domain is :

$$\hat{\alpha} = \frac{\hat{\mathbf{y}}}{\hat{\mathbf{x}} \odot \hat{\mathbf{x}} + \lambda} \quad (6)$$

if we added the Kernels and used the Kernel trick we can show that the dual will have the form

$$\hat{\alpha} = \frac{\hat{\mathbf{y}}}{\hat{\mathbf{k}}^{xx} + \lambda} \quad (7)$$

where $\hat{\mathbf{k}}^{xx}$ is the kernel vector formed from inner product from \mathbf{x} and itself

The KCF take advantage of the fact that the convolution of two patches (loosely, their dot-product at different relative translations) is equivalent to an element-wise product in the Fourier domain. Thus, by formulating their objective in the Fourier domain, they can specify the desired output of a linear classifier for several translations, or image shifts, at once[3]. The power of the kernel trick comes from the implicit use of a high-dimensional feature space , without ever instantiating a vector in that space.[3]

2.2 KCF Family of trackers

One adaption of KCF is SAMF (Scale-Adaptive Kernel Correlation Filter) [4]. This adaption just like KCF , learn the filter and apply it on translated patches , however it searched for different scales and look at the maximum response over all the scales . We are solving the rotation problem it more efficient way (one shot) rather than trying all different rotations and take the one with the highest response.

other versions of KCF are those of deep features that enhance the performance of KCF and allow it to be scale and rotational invariant for enough training of deep Neural Network like [6]. However these requires long training and GPUs and also lack the speed the original KCF has.

2.3 Histogram of oriented gradients (HOG)

A generic way to extract orientation feature of objects is to find the distribution of the gradients in cells that combine a number of pixels. This is very powerful and fast technique to characterize the orientation of an object. If we take the image patch to be one cell and we choose enough number of bins we can have a global descriptor for the patch as can be seen in figure 2, it can be seen that the descriptor has circular structure (coming from the fact that rotating 180 degrees give the same HOG descriptor for the patch) that will prove crucial in formulating the solution for RKCF.

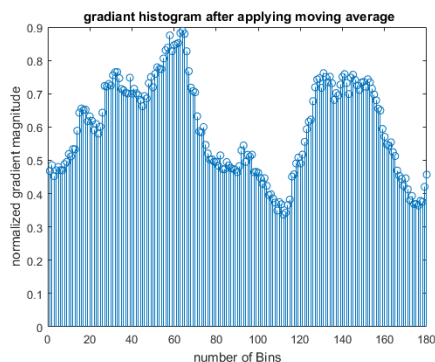


Figure 2: HOG feature of image patch , smoothed for better performance , observe the circular structure of the global descriptor

for an image patch like the one in 3 by multiplying it by a cos window and then rotating it, its global HOG descriptor will suffer a shift like the one in figure 4 which is very similar to what will happen to the target when we learn the rotation filter and apply it on the new rotated patch in RKCF

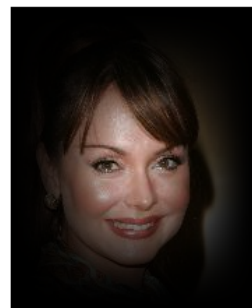


Figure 3: image patch multiplied by cos window, (typical in CF trackers) to reduce the effect of background and the boundary effect

3 methodology

3.1 Derivation of augmented KCF (RKCF)

We can extend the optimizing objective of KCF (or any other tracker that uses Correlation Filters) to include the rotation information of the target as follows

$$\min_{\mathbf{w}, \mathbf{r}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda_1 \|\mathbf{w}\|_2^2 + \|\mathbf{A}\mathbf{r} - \mathbf{g}\|_2^2 + \lambda_2 \|\mathbf{r}\|_2^2 \quad (8)$$

where $\mathbf{w} \in \mathbb{R}^{m \times n}$ is the filter learned in KCF step with window size $m \times n$ and $\mathbf{X} \in \mathbb{R}^{mn \times mn}$ is the data matrix of all possible shifts of the image patch that contains the target being tracked (like before) in which 2D convolution would be performed in the detection phase.

$\mathbf{r} \in \mathbb{R}^b$ is the rotation filter that is to be learned and b is the number of bins in HOG descriptor (a) describing globally the patch , \mathbf{A} is the circulant matrix of vector \mathbf{a} that

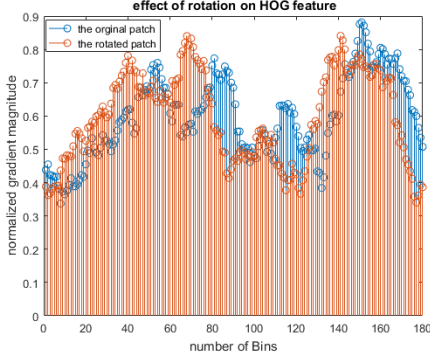


Figure 4: the effect of rotating image like the one in 3 on its global HOG descriptor

reflects all possible rotations of the target, in which 1D convolution would be performed in the detection phase. \mathbf{g} is just like \mathbf{y} of the KCF, it is a 1D typical response of the rotation filter \mathbf{r} on the \mathbf{a} descriptor if there was no rotation from one frame to the other in the detection. λ_1, λ_2 are the regularization variables.

The optimization 7 is separable in \mathbf{w}, \mathbf{r} and has the closed form solutions :

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda_1 \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (9)$$

$$\mathbf{r} = (\mathbf{A}^T \mathbf{A} + \lambda_2 \mathbf{I})^{-1} \mathbf{A}^T \mathbf{g} \quad (10)$$

Following similar procedure as in KCF derivation we can show that \mathbf{r} can be written as :

$$\hat{\mathbf{r}} = \frac{\hat{\mathbf{a}}^* \odot \hat{\mathbf{g}}}{\hat{\mathbf{a}}^* \odot \hat{\mathbf{a}} + \lambda_2} \quad (11)$$

in the dual domain we can formulate the following solution

$$\hat{\alpha}_r = \frac{\hat{\mathbf{g}}}{\hat{\mathbf{k}}^{aa} + \lambda_2} \quad (12)$$

in which $\hat{\mathbf{k}}^{aa}$ is just like $\hat{\mathbf{k}}^{xx}$ in KCF.

3.2 Detection phase in RKCF

the algorithm used in detection in RKCF is similar with KCF in which we apply the learned filter \mathbf{w} to the current

patch and we get response. comparing the position of the maximum of the response to the \mathbf{y} maximum position. This translation of the maximum dictates the translation of the target position in the next frame. In RKCF we do the same thing as KCF then apply the \mathbf{r} filter that we learned in 11 to give a response like the one in the following figure 5. The shift of the maximum of this response compared to the standard response (Gaussian centred at the middle) will give the rotation of the object from one frame to another.

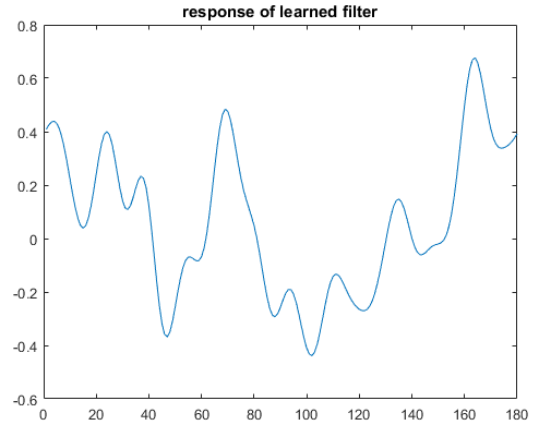


Figure 5: the response of the HOG of the new patch to the rotation filter, max shift from zero represent the rotation in degrees

after we have detected the rotation from one frame to the following one, this angle will be used to enhance the KCF response by counter rotate the new patch by the amount it suffered from the first frame, re apply \mathbf{w} on the "adjusted" patch and see if the max response of that was higher than the original response. If this was the case then take the new response translation as the trusted one and rotate it and apply it on the target. If not, then go with KCF suggested translation. This will insure to some extent that tracker doesn't drift away from the target based on false rotation detection. Algorithm 1 summarizes the RKCF tracking scheme.

each frame, do the following:

- learn any CF filter , like in equation 5
- apply what you learn on the new frame patch and record max response U
- learn filter r based on HOG of the old patch like in 11
- apply r on the new patch HOG
- get the rotation θ as described in section 3
- rotate the patch to $-\theta$ and apply the original CF on it **if** $\max(\text{response}) > U$ **then**
 - trust the rotation and take its max translation ;
 - rotate the translation by θ and use it as target translation ;
- else**
 - use whatever U gave as translation of the target;

Algorithm 1: RKCF learning and detection algorithm

4 Experiments

4.1 large scale experiment for rotation detection

A set of 816 images were each rotated by 1000 random rotations to give a validation set of 816000 samples . To test the rotation filter effectiveness in detecting rotation, two other ways of detecting rotation based on HOG feature were tested and bench-marked. the filter have access on the upright patch and the rotated patch but no access to the actual angle at which it was rotated and its goal is to find that angle. The following results were obtained for a cosine window and a Gaussian window on the patch to cancel the boundary effect. The two other ways are correlation between the two HOGs , and observing the shift of the max of the HOG descriptor from one frame to another. The result is mean abs error in degrees for all the permutations.

4.2 The RKCF on OBT50 .

OBT50 is one of the most famous data-sets for visual tracking since it was released in 2013 and extended to OBT100 in 2015 [7] .We assess our RKCF algorithm on OBT50 dataset and compare to the base line (KCF in this case).We observe huge enhancement of the base line for some difficult videos like the "Matrix" video on

boundary window	cos window	Gaussian window
rotation filter	15.29	16
correlation	13.46	19.04
max shift	36.96	43.19

Table 1: The absolute error in degrees for different rotation detection envelopes and using different rotation detection techniques based on HOG features

which there is a lot of rotation that a regular KCF suffers dramatically. figure 6 show how the target suffers huge rotation from a frame to the following frame.



Figure 6: The "Matrix" sequence in which in one frame the head is straight and five frames later its 90 degrees rotated making it difficult to track by KCF

In this specific sequence the difference in precision is 20 pts more for RKCF !. The following figure 7 depicts this.

We can see in this specific sequence that the rotation from one frame to another is huge as shown in figure 8

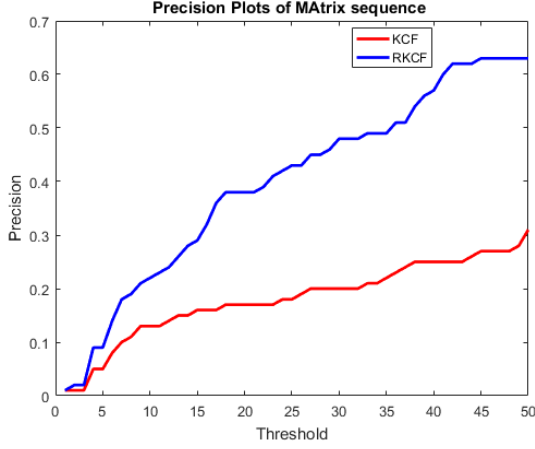


Figure 7: The "Matrix" sequence precision plot comparing baseline KCF and proposed RKCF

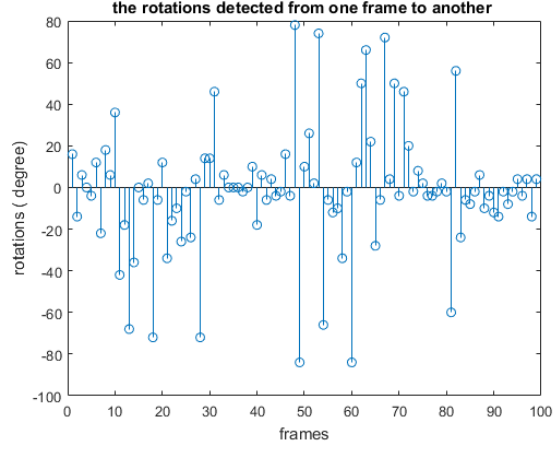


Figure 8: The "Matrix" sequence detected target rotation from one frame to another

. So we propose a new method to evaluate the rotational difficulty of a video directly from its KCF response sequence. We call it \mathcal{U} (pronounced moh) the rate of rotational change in target in which

$$\mathcal{U} = \text{std}(\theta_i) \forall i \in n \quad (13)$$

where n is the number of frames in the sequence and θ_i is the rotation detected by RKCF in the frame i . in the matrix sequence \mathcal{U} was 24.8, very high rate per frame

We define a success rate R for our proposed RKCF that assess its quality on a video as follows :

$$R = s/(s + f) \quad (14)$$

in which s is the number of frames in which proposed RKCF scheme gave higher response than baseline KCF, f is the number of frames in which proposed RKCF scheme gave lower response than baseline KCF.

Performing the test on the whole data-set, we get the following analysis results 910 showing rotational difficulty \mathcal{U} , success rate R for all videos.

A mean success rate of all data-set of 31.57% was obtained. The following precision plot compares RKCF and baseline KCF on the whole dataset.

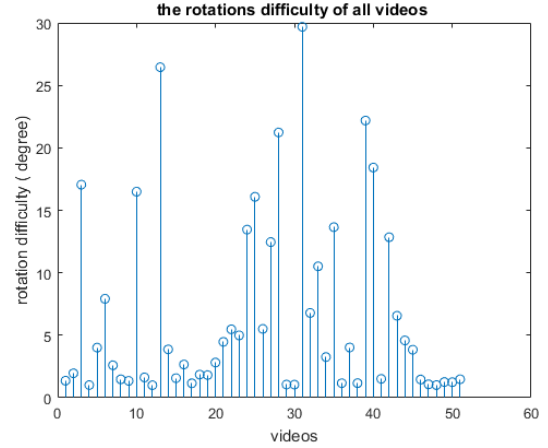


Figure 9: The rotational difficulty \mathcal{U} on the OBT50 data set

5 Observations

- Importance of matching target size(by scale of filter) on over-all tracking precision, correct rotation with bad scale doesn't help that much
- For constant scale targets, proposed RKCF achieve as good or better than KCF due to its rotational capability

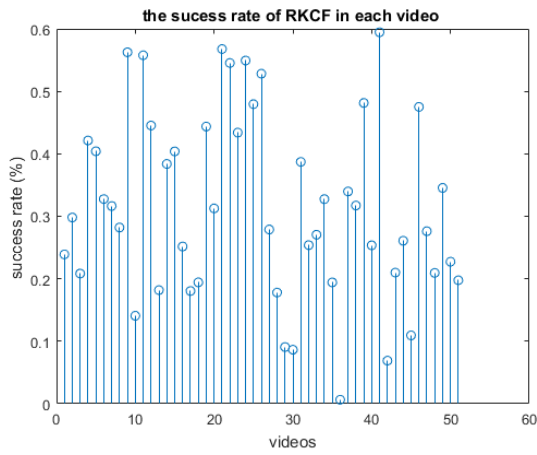


Figure 10: The success rate of RKCF R on the OBT50 data set

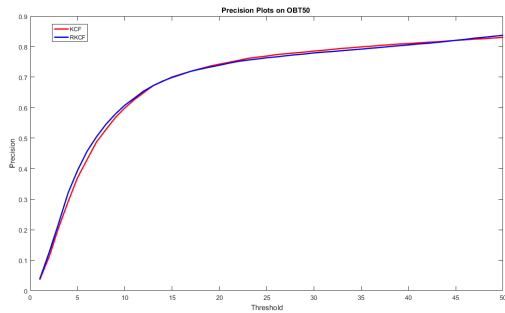


Figure 11: The precision plot of RKCF and baseline KCF on the OBT50 data set

- videos with \bar{U} around 25 , gives max performance of RKCF over KCF.
- it might be easier and better to apply correlation between HOGs instead of learning rotational filter r .

6 Future work

- we propose generalizing this framework to include all types of features (not only HOG) so RKCF can be used

to augment Deep Features trackers like one in [6], the generalization can be

- we propose justification on why the typical target response of the rotational filter to be Gaussian centered at the middle .
- We propose updating the learning rule of w to include the rotation information θ .
- We propose a KLT framework to tackle scale and rotation in Correlation Filter fashion , by linearizing around the current frame .

References

- [1] A. Bibi. Target response adaptation for correlation filter tracking, 2010. in: European Conference on Computer Vision, 2015).
- [2] M. Danelljan. Accurate scale estimation for robust visual tracking, 2014. roceedings of the British Machine Vision Conference. BMVA Press, September 2014.
- [3] J. Henriques. High-speed tracking with kernelized correlation filters, 2015. in: IEEE Transactions on Pattern Analysis and Machine Intelligence (Volume: 37, Issue: 3, March 1 2015).
- [4] Y. Li and J. Zhu. A scale adaptive kernel correlation filter tracker with feature integration, 2014. in: European Conference on Computer Vision, 2014.
- [5] D. Marcos. Learning rotation invariant convolutional filters for texture classification, 2016. arXiv.
- [6] J. Valmadre. End-to-end representation learning for correlation filter based tracking, 2017. arXiv.
- [7] Wu. Online object tracking: A benchmark, 2013. Proceedings of the IEEE conference on , 2013 - cv-foundation.org.