



Object Tracking with Kernelized Correlation Filters (KCF)

--CS585 Final Project Report

Nan Zhou
Lingzhuo Zhao
Yanjie Chen

Abstract

This project aims to study a traditional method, Kernelized Correlation Filters(KCF) to achieve high-speed tracking. KCF is an improvement for the circulant structure tracker (CST) [1] by adding in the Histogram of Oriented Gradients (HOG) feature. We aim to apply the KCF tracker to some of the benchmark datasets and analyze the performance in comparison with other tracking algorithms.

Given a series of video frames, our objective is to track at least one target in the object. It involves researching and developing an algorithm with good performance, testing multiple benchmark datasets and performing precise evaluation.

Related Work

Visual tracking is a popular research problem in computer vision due to the applications in different areas such as activity recognition, human-computer interaction, and motion detection. Numerous tracking algorithms have been proposed which allow for some assumptions about the target object. It is ideal to track an object with little information about it. A very successful method is tracking-by-detection.

Most tracking algorithms can be categorized into two types: generative and discriminative based on their appearance models. A generative tracking method learns an appearance model to represent the target and searches for image regions with the highest matching scores as the results[6]. A discriminative classifier is the key component of modern trackers in general. It aims to distinguish the target and its surroundings. Some examples of the related tackers are Support Vector machines(SVM), Random Forest Classifiers, boosting variants, and Naïve Bayes classifiers proposed by Zhang et al. [3], which can predict the target's location directly rather than its presence in an image patch. This paper will focus on two discriminative methods. One is proposed by Kalal et al[4], and the other is proposed by Bolme et al. [5].

The approach proposed by Kalal et al[4] can cut down the number of training samples as it uses a list of structural constraints to lead the sampling process of a boosting classifier. However, there

exist some limitations such as a highly time-consuming training process, limited training features that can be used and tedious structural heuristics tuning. The correlation filters [5] only use partial computation power although it is more sophisticated. The trick is the convolution of two patches is the same as an element-wise product in the Fourier domain, which makes the process of translation and image shifts by linear classifier more efficient. But, it can be limiting as well.

Method

KCF is a more representative discriminant tracking algorithm, which is tasked to distinguish between the target and the surrounding environment. On the basis of the original CSK algorithm, the main feature is to use the circulant matrix Fourier space diagonalization to simplify the calculation. Before learning KCF, you can read the other two papers, namely the 10-year MOSSE tracking algorithm and the CSK algorithm. The MOSSE algorithm applied convolution to the tracking algorithm for the first time, which greatly improved the algorithm speed. Later, CSK on this basis, the algorithm adds a regular term to avoid over-fitting, and at the same time introduces a circulant matrix and a kernel function to improve computational efficiency, and the performance is greatly improved. Figure 1 includes some qualitative results for KCF as compared to top-performing algorithms. A 'x' is used to represent missing trackers.

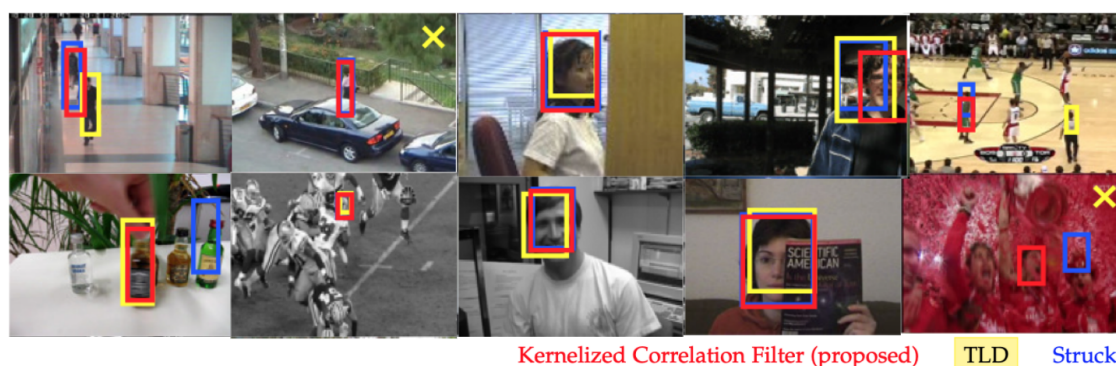


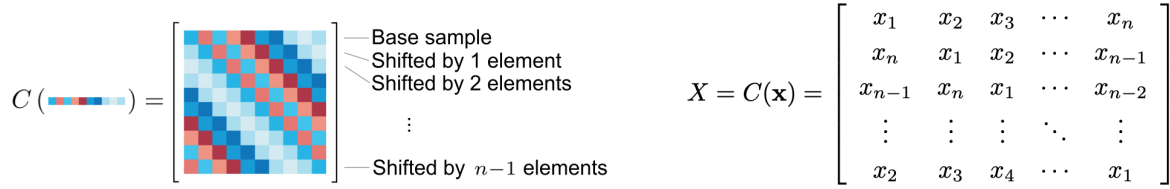
Figure 1. Qualitative results for KCF, compared with Struck and TLD.

Figure 2 shows the experimental results of the benchmark datasets using different algorithms. The KCF algorithm with HOG features included is able to give better results as compared to other methods.

Algorithm	Feature	Mean precision (20 px)
KCF	Raw pixels	56.0%
KCF	HOG	73.2%
Struck		65.6%
TLD		60.8%

Figure 2. Summary of experimental results.

KCF is a tracking method of Tracking-By-Detection, which is the same as TLD and OAB. The tracking object is represented as a positive sample, while the surrounding environment is negative to train a discriminative classifier. This method uses the circulant matrix of the expanded area of the target to collect positive and negative samples. Figure 1 is an illustration of a circulant matrix.



$$C(\text{vector}) = \begin{bmatrix} \text{Base sample} \\ \text{Shifted by 1 element} \\ \text{Shifted by 2 elements} \\ \vdots \\ \text{Shifted by } n-1 \text{ elements} \end{bmatrix}$$

$$X = C(\mathbf{x}) = \begin{bmatrix} x_1 & x_2 & x_3 & \cdots & x_n \\ x_n & x_1 & x_2 & \cdots & x_{n-1} \\ x_{n-1} & x_n & x_1 & \cdots & x_{n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_2 & x_3 & x_4 & \cdots & x_1 \end{bmatrix}$$

Figure 3. Circulant matrix. The rows are cyclic shifts of a vector image.

Consider a $n \times 1$ vector, denoted as x . The circulant matrix is achieved by a cyclic shift operator, which is the permutation matrix below.

$$P = \begin{bmatrix} 0 & 0 & 0 & \cdots & 1 \\ 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix}$$

Next, by applying ridge regression to train the detector, we successfully use the diagonalization property of the circulant matrix in Fourier space to convert complex matrix operations into vector dot multiplication.

$$\hat{\mathbf{w}} = \frac{\hat{\mathbf{x}}^* \odot \hat{\mathbf{y}}}{\hat{\mathbf{x}}^* \odot \hat{\mathbf{x}} + \lambda}$$

The fraction denotes element-wise division. This transformation greatly reduces the amount of calculations and improves the efficiency of tracking. The ridge regression of the linear space is mapped to the higher-dimensional nonlinear space through the kernel trick, which consists of:

- 1) Express \mathbf{w} as a linear combination of samples:

$$\mathbf{w} = \sum_i \alpha_i \varphi(\mathbf{x}_i)$$

- 2) Write the algorithm in dot-products

$$\varphi^T(\mathbf{x})\varphi(\mathbf{x}') = \kappa(\mathbf{x}, \mathbf{x}')$$

The original problem thus becomes non-linear. Alpha can then be computed efficiently by,

$$\hat{\alpha} = \frac{\hat{\mathbf{y}}}{\hat{\mathbf{k}}^{\mathbf{x}\mathbf{x}} + \lambda}$$

together with the fast detection response,

$$\hat{\mathbf{f}}(\mathbf{z}) = \hat{\mathbf{k}}^{\mathbf{x}\mathbf{z}} \odot \hat{\alpha}$$

Gaussian kernel with $\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2\right)$ is a special but useful one, and we will get

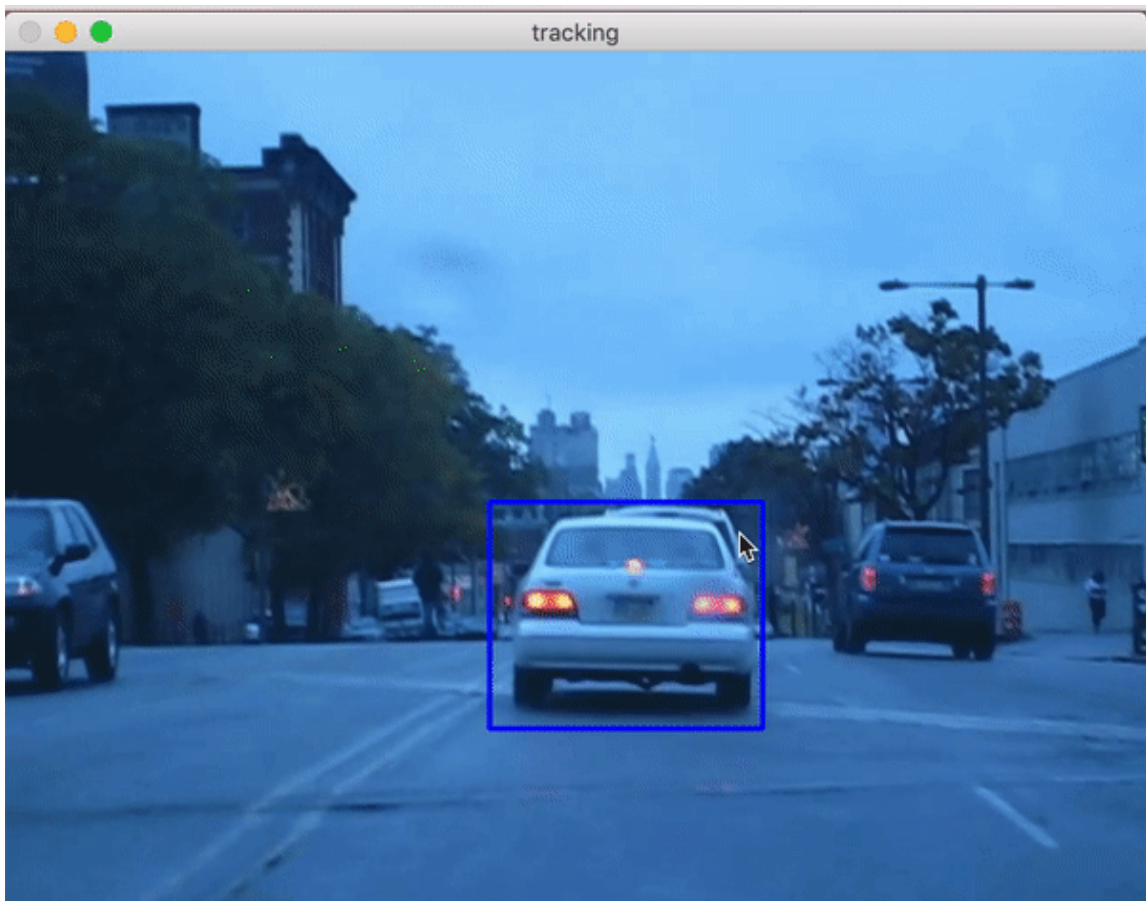
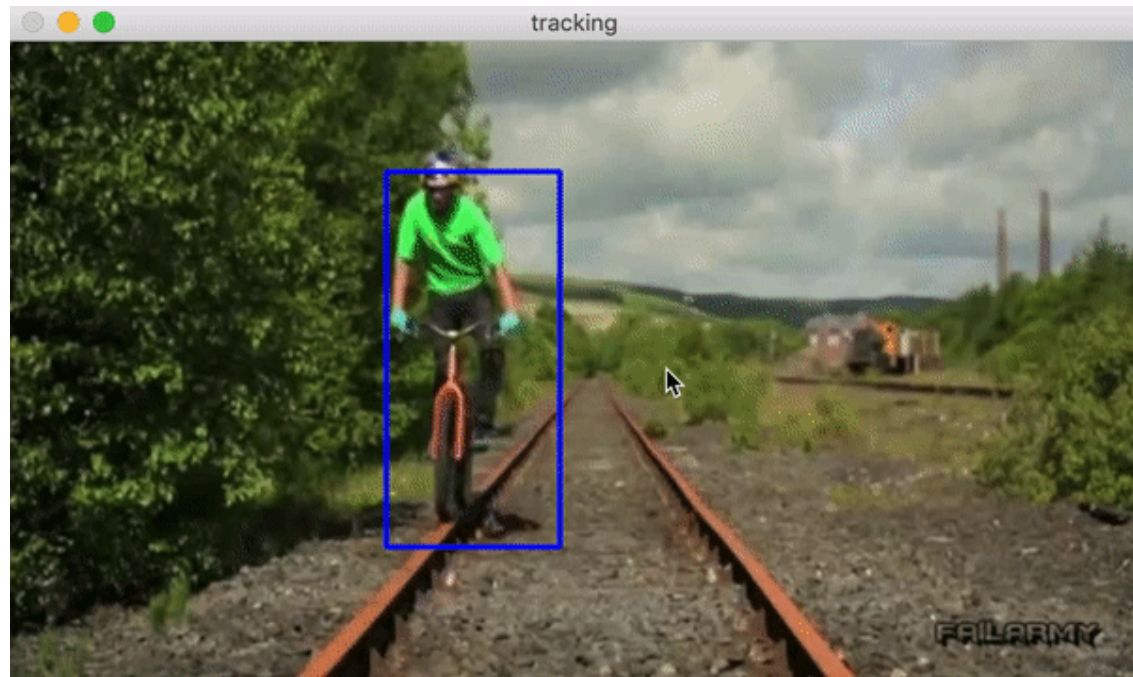
$$\mathbf{k}^{\mathbf{x}\mathbf{x}'} = \exp\left(-\frac{1}{\sigma^2} \left(\|\mathbf{x}\|^2 + \|\mathbf{x}'\|^2 - 2\mathcal{F}^{-1}(\hat{\mathbf{x}}^* \odot \hat{\mathbf{x}}')\right)\right)$$

In the nonlinear space, by solving a dual problem and some common constraints, it is also possible to use the circulant matrix Fourier space diagonalization to simplify the operation. Expose the limitation that can only handle single-channel features, and give a way to deal with multi-channel features. For multi-channel, it is assumed that a vector \mathbf{x} concatenates the individual vectors for C channels. Below is a concrete example applied on the Gaussian kernel. The final result is the summation of each channel in the Fourier domain due to linearity of the Discrete Fourier Transform.

$$\mathbf{k}^{\mathbf{x}\mathbf{x}'} = \exp\left(-\frac{1}{\sigma^2} \left(\|\mathbf{x}\|^2 + \|\mathbf{x}'\|^2 - 2\mathcal{F}^{-1}(\sum_c \hat{\mathbf{x}}_c^* \odot \hat{\mathbf{x}}'_c)\right)\right)$$

In general, in this algorithm the target object or region of interest (ROI) is chosen in the first frame. The next step is to train a model with the image patch, which is larger than the ROI, at the initial position. The purpose is to provide some context. In the following new frames, it is able to detect over the patch at the previous position. And the target position is updated to the one that yielded the maximum value accordingly. Then a new model is trained at the updated position.

Results



Fail to track



Discussion and Conclusion:

Our project took the advantage of our improved KCF algorithm and achieved relatively satisfactory results that managed to capture moving objects. While we were researching object tracking related papers, we found out that KCF would, in theory, capture the boxed object in multiple scenarios including rapid moving objects, moving objects and zooming cameras etc,. On the other hand, there are still limitations due to the principle of KCF model.

The results heavily depend on the given training dataset and boxed frame size. We noticed the target object size changed with time and the boxed frame drifted in certain scenarios. For instance, in the third demo, we can see that the camera was moving up and down while also zoomed in and out at the same time. In this case, the KCF algorithm was not able to capture the target in real time. The reason behind this could be that the depth of the feature vectors failed to capture the video information, occlusion occurred, or low resolution frame loss, etc,.

Future Work

There is still some future work that could be done to improve our KCF algorithm. First, we could test our algorithm in some different datasets. Based on the results of different datasets, some parameters in our algorithm would be tuned and we can get better results. Second, right now we need to manually select an object to track using a bounding box at the initial frame for our algorithm. We can improve this by making the program automatically tracking the object in the video. Also, we could develop multi-object tracking based on this algorithm.

References

1. Henriques, J., Caseiro, R., Martins, P., Batista, J.: Exploiting the circulant structure of tracking-by-detection with kernels. In: ECCV. pp. 702-715 (2012)
2. J. F. Henriques, R. Caseiro, P. Martins and J. Batista, "High-Speed Tracking with Kernelized Correlation Filters," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 3, pp. 583-596, 1 March 2015, doi: 10.1109/TPAMI.2014.2345390.
3. K. Zhang, L. Zhang, and M.-H. Yang, "Real-time compressive tracking," in *ECCV*, 2012.
4. Z. Kalal, K. Mikolajczyk, and J. Matas, "Tracking-learning- detection," *TPAMI*, 2012.
5. D. S. Bolme, J. R. Beveridge, B. A. Draper, and Y. M. Lui, "Visual object tracking using adaptive correlation filters," in *CVPR*, 2010, pp. 2544-2550.
6. K. Zhang, L. Zhang, Q. Liu, D. Zhang, and M.-H. Yang. Fast visual tracking via dense spatio-temporal context learning. In *Proceedings of the European Conference on Computer Vision*, 2014.