

The Large margin tracking method with circulant feature maps

创建时间：2017/9/23 18:37

作者：XY.Wang

the large margin tracking method with circulant featur

利用循环移位样本的特征图(circulant feature maps)与一个向量的内积，作为是目标的可能性；

目标函数的约束条件：使得目标中心处的patch是目标的可能性大于 (w,h) 处是目标的可能性超过一个距离。意为large margin！！

output pairs (x, y) with $F : X \times Y \rightarrow \mathbb{R}$ from which we can acquire a prediction by maximizing F over the response variable for a specific given input x . Then the general form of the function f can be denoted as

$$f(x; w) = \arg \max_{y \in Y} F(x, y; w) \quad (1)$$

衡量是目标的可能性

where we assume F to be a linear function, $F(x, y; w) = \langle w, \Psi(x, y) \rangle$ and w denotes the parameter vector which can be learned from the soft-margin support vector machine learning over structured outputs. F can also be extended to nonlinear situation which will be discussed in the next section. We penalize margin violations by a quadratic term, leading to the following optimization problem:

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_{w=1}^{W-1} \sum_{h=1}^{H-1} \xi_{w,h}^2$$

$$\text{s.t. } \forall w, \forall h, \forall y_{w,h} \in Y \setminus y_{0,0} :$$

$$F(x, y_{0,0}; w) - F(x, y_{w,h}; w) \geq \sqrt{\Delta(y_{0,0}, y_{w,h})} - \xi_{w,h}$$

松弛因子

意义：寻找 w 的解，使得两个patches之间是目标的可能性的差距 \geq 二者的高斯权重差的where $y_{0,0}$ denotes the observed output with no cyclic transform and $\xi_{w,h}$ is the slack variable which penalizes the margin violations. The regularization parameter $C > 0$ controls the trade-off between training error minimization and margin maximization. $\Delta(y_{0,0}, y_{w,h})$ quantifies the loss associated with a prediction $y_{w,h}$ when the true output value is $y_{0,0}$. We define the loss function as

$$\Delta(y_{0,0}, y_{w,h}) = m(y_{0,0}) - m(y_{w,h}) \quad (3)$$

a fast optimization algorithm that builds up a bridge b problem formulation and the well-known correlation fi

给定了需要优化的目标函数和约束条件，利用一种与CF结合的快速在线优化算法，将求解过程转化到频域，加速了优化过程。

w 和 z 其中一个固定的话，优化过程就变得很简单了，一步即可求出来，其中 w 是在频域计算的。然后交替更新，优化。

(1) w 固定， z 的求解方式：

There are two variables w and z to be solved in Eq.7. Whenever one of them is known, the subproblem on the other has a closed form solution. Thus similar to [34], we introduce the alternating optimization algorithm to solve the model efficiently by iterating between the following two steps.

Update z . Given w , the subproblem on z becomes:

$$\min_z \|z - (w^T \Phi_0 - w^T \Phi - \Upsilon)\|_2^2, \text{ s.t. } z \geq 0 \quad (9)$$

Then the closed form solution of z is:

$$z = \max \{w^T \Phi_0 - w^T \Phi - \Upsilon, 0\} \quad (10)$$

w 和 z 其中一个固定的话，优化过程就变得很简单了，一步即可求出来。然后交替更新，优化。

(2) z 固定， w 转化到频域进行求解，加速了运算（计算公式如下图）

(3) 非线性模型的求解：是需要求解参数 α ，由公式可计算出 w （计算公式如下图）

Update w . Given z , the subproblem on w becomes:

$$\min_w \frac{1}{2} \|w\|^2 + C \|w^T \Phi - (w^T \Phi_0 - \Upsilon - z)\|_2^2 \quad (11)$$

哪里有用到CF理论？

In order to employ the correlation filter theory, we define $u_0 = w^T \Phi_0$ which stands for a plane whose height is the highest peak of $w^T \Phi$ in the last iteration. Then the closed form solution of w is:

将循环样本的计算转换到了频域求解？
就是滤波器理论？

$$\hat{w} = \frac{\hat{\Psi}^*(x, y_0) \circ \hat{u}^T}{\hat{\Psi}^*(x, y_0) \circ \hat{\Psi}(x, y_0) + \frac{1}{2C}} \quad (12)$$

where $u = u_0 - \Upsilon - z$ and \circ denotes the element-wise division.

Nonlinear extension. The proposed linear model can be extended to a nonlinear model by the kernel trick $K_{ij} = \langle \varphi(\Psi(x, y_i)), \varphi(\Psi(x, y_j)) \rangle$ where $\varphi(\bullet)$ indicates the implicit use of a high-dimensional feature space. The solution w can be represented as $w = \sum_{w=0}^{W-1} \alpha_w \varphi(\Psi(x, y_w))$.

The optimization now is rewritten as

$$\begin{aligned} \min_{\alpha} & \alpha^T \mathcal{F}^{-1} (\hat{K}^{\Psi_0 \Psi_0} \circ \hat{\alpha}) \\ & + C \left\| \mathcal{F}^{-1} (\hat{K}^{\Psi_0 \Psi_0} \circ \hat{\alpha}) - (u_0 - \Upsilon - z)^T \right\|_2^2 \quad (13) \\ \text{s.t. } & z \geq 0 \end{aligned}$$

where $\Psi_0 = \Psi(x, y_0)$ and $\hat{K}^{\Psi_0 \Psi_0}$ denotes the DFT of the first row of the circulant kernel matrix K whose elements are K_{ij} . The closed form of the subproblem on α is

$$\hat{\alpha} = \frac{\hat{u}^T}{\hat{K}^{\Psi_0 \Psi_0} + \frac{1}{2C}} \quad (14)$$

where \circ denotes the element-wise division.

多峰目标检测

新的一帧到来时，通过公式(1)向量参数与特征图做内积得到的响应图来预测目标位置，

tion \hat{f} can be denoted as

$$f(x; w) = \arg \max_{y \in Y} F(x, y; w) \quad (1)$$

衡量是目标的可能性

x 是以上一帧目标位置为中心，与目标区域成比例的一个image patch；

是目标可能性的函数 F 定义为两个向量的内积：

$$F(x, y; w) = \langle w, \Psi(x, y) \rangle$$

由前面的内容可知，求解参数 w 的过程可以通过FFT来加速。

单峰目标检测(只找最高峰)

所有的循环样本的响应图放在一起，就变成了：

$$\text{单峰} \quad F(s, y; w) = \mathcal{F}^{-1} \left(\hat{\Psi}_{s0}^* \circ \hat{w} \right) = \mathcal{F}^{-1} \left(\hat{k}^{\Psi_{x0} \Psi_{s0}} \circ \hat{a} \right) \quad (15)$$

寻找响应图的最高峰，作为预测的目标位置。（s是前面的x）

解释：s的循环移位样本的特征与w做内积运算==原始image patch s的联合特征图与w做循环卷积（移位量为0时）->转换到频域，就是二者的element-wise product

缺陷：然而这种检测方法是单峰检测，如果有相似物体或背景干扰存在的话，它们对应的响应值可能很接近目标甚至比目标更高，因而响应图中的最高峰可能是把最高峰当做目标，而是将所有的峰都考虑进去（即多峰检测），从中寻找代表目标位置的那个峰。

多峰目标检测（从多个峰里再找最高分）

在响应图的基础上，乘上一个与响应图大小相同的0/1矩阵B（B的元素在局部最大值处的值设为1，其余为0）。则P(s)中的非零值就表示响应图中的多峰。——其

文中这样描述：

Consequently, a multimodal target detection method is proposed to improve localization precision further. For the unimodal detection response map $F(s, y; w)$, the multiple peaks are computed by

$$\text{多峰} \quad P(s) = F(s, y; w) \circ B \quad (16)$$

where B is a binary matrix with the same size as $F(s, y; w)$, which identifies the locations of local maxima in $F(s, y; w)$. The elements at the locations of local maxima in B are set to 1, while others are set to 0. All non-zero elements in $P(s)$ indicate multiple peaks in the response map of s. P(s)只保留了响应图中的局部最大值（多峰）。

多峰响应公式：

$$\text{多峰} \quad P(s) = F(s, y; w) \circ B \quad (16)$$

多峰检测方法：当多峰和最高峰的比例超过预设的阈值时，就以这些峰所在的位置为中心，分别提取patch作为s，然后用公式15寻找每个峰所在区域的最大值。下图中间是检测到的多峰，最右是多峰检测后的最终目标位置。

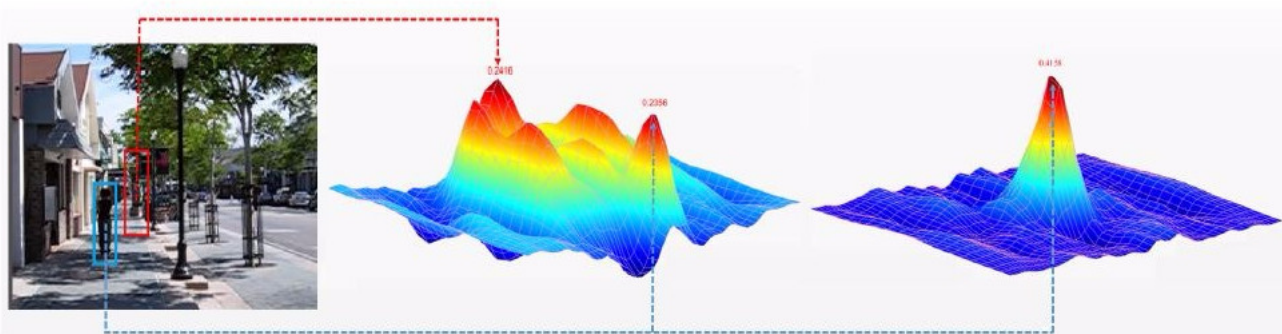


Figure 1. Illustration of multimodal target detection in sequence *human9* from OTB-15 [30]. The blue bounding box indicates the correct location of target, the red one is an incorrect detection. The response of the target is weaker than the background area within the red bounding box as shown in the middle. The unimodal detection will regard the highest peak as the target leading to false detection. The proposed multimodal target detection will redetect the areas centered at other peaks to find the maximum peak among these response maps as the right subfigure and locate the correct position of the target.

High-coinfidence update

大多数的跟踪器不考虑检测结果是否精确，简单粗暴地每帧都更新。实际上，当目标被遮挡或消失了，或者当前帧的检测结果根本就不准确时，再用这些结果去更

下图展示了目标被遮挡后继续每帧更新模型，造成后续某一帧跟踪结果错误的现象。（第一列的绿色框对应每帧更新，空色是high-coinfidence更新；第二列是更新的跟踪结果，可以看到(f)图跟踪错了）

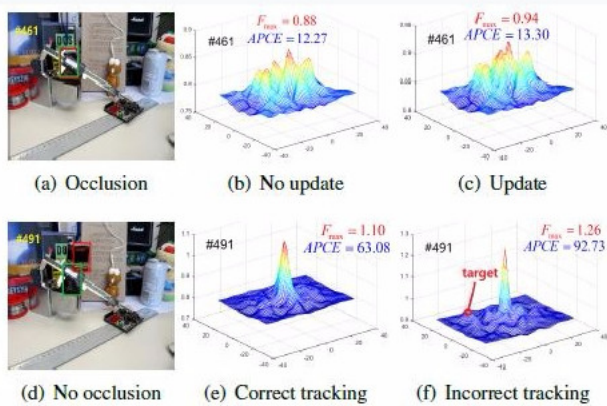


Figure 2. The first column are the shots of sequence *box* from OTB-15, where the red bounding boxes indicate the tracking results of LMCF with high-confidence update strategy and the green ones belong to the LMCF-NU which updates the tracking model in each frame. The response maps in the second column are corresponding to LMCF and the third column corresponding to LMCF-NU. The red annotation in the last subfigure points out the right position of the target in this response map.

本文：利用跟踪结果的反馈信息（跟踪结果是否足够可信），来决定是否有必要更新模型！

理想的响应图：只有一个sharp peak（锐利的峰），其他区域平滑下降。相关峰越尖锐，检测结果越精确。

非理想的响应图：有很多的波动。如果继续用不确定是否是目标的区域作为训练样本去更新模型，可能导致跟踪结果出错。

A. 定义两个条件：

$F_{\max} = \max F(s, y; w)$ ：响应图中的最高相应分数

APCE：表示响应图的波动程度和检测结果的可信度。对于更尖锐的峰和少的噪音，响应图只有一个尖锐的峰和平滑下降的区域，APCE变得更大。当目标被

confidence feedback mechanism with two criteria. The first one is the maximum response score F_{\max} of the response map $F(s, y; w)$ defined as

$$F_{\max} = \max F(s, y; w) \quad (17)$$

The second one is a novel criterion called average peak-to-correlation energy (APCE) measure which is defined as

$$APCE = \frac{|F_{\max} - F_{\min}|^2}{\text{mean} \left(\sum_{w,h} (F_{w,h} - F_{\min})^2 \right)} \quad (18)$$

where F_{\max} , F_{\min} and $F_{w,h}$ denote the maximum, minimum and the w -th row h -th column elements of $F(s, y; w)$. APCE indicates the fluctuated degree of response maps and the confidence level of the detected targets. For sharper peaks and fewer noise, i.e., the target apparently appearing in the detection scope, APCE will become larger and the response map will become smooth except for only one sharp peak. Otherwise, APCE will significantly decrease if the object is occluded or missing. 更新条件：

B. 更新条件(High-coinfidence)：

当当前帧的 F_{\max} 和APCE与它们各自的历史均值的比值超过 β_1 和 β_2 时，就认为当前帧的跟踪结果是高可信的，这个时候才选择更新模型。
the object is occluded or missing. 更新条件:

When these two criteria F_{\max} and APCE of the current frame are greater than their respective historical average values with certain ratios β_1, β_2 , the tracking result in the current frame is considered to be high-confidence. Then the proposed tracking model will be updated online with a learning rate parameter η as

$$\begin{aligned}\hat{\alpha}^t &= (1 - \eta) \hat{\alpha}^{t-1} + \eta \hat{\alpha} \\ \hat{\Psi}_{x_0}^t &= (1 - \eta) \hat{\Psi}_{x_0}^{t-1} + \eta \hat{\Psi}_{x_0}\end{aligned}\tag{19}$$

C. 这样做的结果是

防止跟踪结果不准确时还去更新模型，造成模型堕落的问题。

本文的算法流程图

Algorithm 1 LMCF tracking algorithm

Input: Frames $\{I_t\}_1^T$, initial target location p_1 , $z = 0$, $u_0 = ones(W, H)$

Output: Target locations of each frame $\{p_t\}_2^T$.

1: repeat

2: Crop an image region s from I_t at the last location p_{t-1} and extract its joint feature map $\Psi(s, y_{0,0})$.

3: Detect the target location p_t with the multimodal detection via Eq.15 and Eq.16.

4: Estimate the scale of the target as [5].

5: Calculate F_{\max} and APCE with Eq.17 and Eq.18.

6: if F_{\max} and APCE satisfy the update condition, then

7: Train the u_0 , z and $\hat{w}(\hat{\alpha})$ with Eq.10 and Eq.12 (14).

8: Update the tracking model with Eq.19.

9: Update the scale estimation model as [5] with η .

10: end if

11: until end of video sequence.

检测阶段

决定是否更新
The la
to extr
layer i
Our tr
a PC
k40 GI

性能分析

只和相关滤波器的跟踪算法进行了比较

1.使用传统特征，fps可以达到85帧每秒；使用CNN特征，明显降低，变化成8f/s。但使用高层特征的精度更高。

使用传统特征，fps可以达到85帧每秒；使用CNN特征，明显降低，变成8f/s。但使用高层特征的精度更高。

Table 2. Characteristics and tracking results of LMCF, DeepLMCF, LMCF-Uni, LMCF-NU and LMCF-N2. The entries in red denote the best results and the ones in blue indicate the second best.

Trackers	multimodal detection	high-confidence update	feature representations	OPE		TRE		SRE		mean FPS
				precision	success	precision	success	precision	success	
LMCF-N2	No	No	conventional	0.799	0.586	0.813	0.612	0.740	0.540	60.74
LMCF-Uni	No	Yes	conventional	0.809	0.606	0.815	0.616	0.757	0.549	61.38
LMCF-NU	Yes	No	conventional	0.813	0.605	0.820	0.619	0.750	0.545	46.45
LMCF	Yes	Yes	conventional	0.839	0.624	0.829	0.625	0.760	0.552	85.23
DeepLMCF	Yes	Yes	deep CNNs	0.892	0.643	0.877	0.649	0.850	0.596	8.11

2.对于遮挡、光线变化等各种情况的鲁棒性都表现的最好：(平均成功率最高)

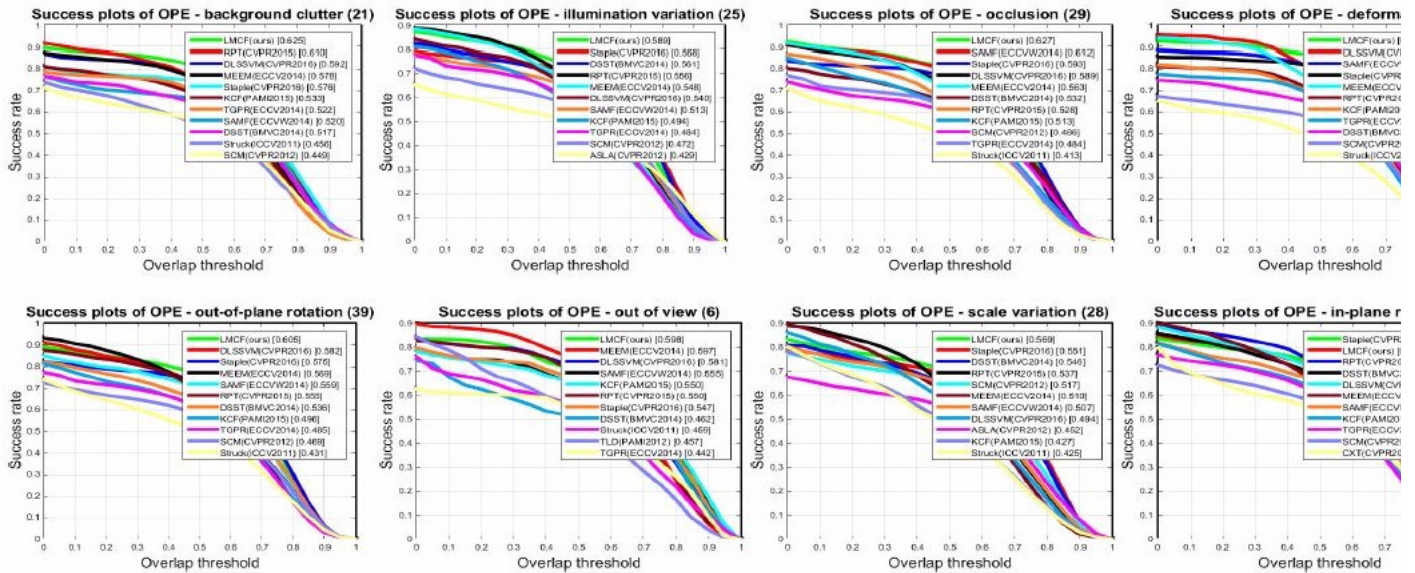


Figure 4. The success plots for 8 challenging attributes including background clutter, illumination variation, occlusion, deformation, out-of-plane rotation, out-of-view, scale variation and in-plane rotation. The proposed LMCF performs best in almost all the attributes are best viewed on high-resolution displays.

3. DeepLMCF的精度仅次于C-COT

