

DSST & FDSST

1.MOSSE

- 两个信号的相关性表示：

$$\begin{aligned}\text{连续}(f \otimes g)(\tau) &= \int_{-\infty}^{\infty} f^*(t)g(t + \tau)dt \\ \text{离散}(f \otimes g)(n) &= \sum_{m=-\infty}^{\infty} f^*[m]g(m + n)\end{aligned}$$

f 表示输入的图像特征， h 表示滤波器模板， g 表示响应输出。

- FFT

$$F(g) = F(f \otimes g) = F(f) \cdot F(h)^*$$

$$G = F \cdot H^*$$

$$H^* = \frac{G}{F}$$

- 构建目标函数

$$\min_{H^*} = \sum_{i=1}^m |H^* F_i - G_i|^2$$

上式的操作都是元素级别的，因此想要得到 H ，只要使其中的每个元素的MOSSE最小即可。上式转化为：

$$\min_{H_{wv}^*} \sum_{i=1}^m |H_{wv}^* F_{wvi} - G_{wvi}|^2$$

- 求偏导，使其为0，得到结果：

$$H_{wv} = \frac{\sum_i F_{i wv} G_{i wv}^*}{\sum_i F_{i wv} F_{i wv}^*}$$

- 最终得到 H

$$H = \frac{\sum_i F_i \cdot G_i^*}{\sum_i F_i \cdot F_i^*}$$

- F_i 与 G_i 的获取

对跟踪框进行随机仿射变换，获得一系列的训练样本 f_i ，而 g_i 则是由高斯函数产生，并且其峰值位置是在 f_i 的中心位置。

- 更新策略

$$\begin{aligned}H_t &= \frac{A_t}{B_t} \\ A_t &= \eta F_t \cdot G_t^* + (1 - \eta) A_{t-1} \\ B_t &= \eta F_t \cdot F_t^* + (1 - \eta) B_{t-1}\end{aligned}$$

2. DCF for Multidimensional Features

- 考虑到图像的特征是多维度的，设 f 为特征，有 d 维， f^l 为其中的第 l 维，损失函数变为：

$$\varepsilon = \left\| \sum_{l=1}^d h^l \star f^l - g \right\|^2 + \lambda \sum_{l=1}^d \|h^l\|^2$$

- H^l

$$H^l = \frac{\bar{G}F^l}{\sum_{k=1}^d \bar{F}^k F^k + \lambda}$$

- 更新策略

$$A_t^l = (1 - \eta)A_{t-1}^l + \eta \bar{G}_t F_t^l$$

$$B_t = (1 - \eta)B_{t-1} + \eta \sum_{k=1}^d \bar{F}_t^k F_t^k$$

- 计算新一帧图像的响应

$$y = \mathcal{F}^{-1} \left\{ \frac{\sum_{l=1}^d \bar{A}^l Z^l}{B + \lambda} \right\}$$

- y 的最大值被认为是目标新位置的估计

3. Fast Scale Space Tracking

- 将位置滤波器与尺度滤波器分开计算，利用二维的位置滤波器获得目标的中心位置后，再用一维的尺度滤波器估计目标在图片中的尺度。
- 训练样本 f 从目标中心扣取，假设当前帧的目标大小为 $P \times R$ ，尺度为 S ，我们扣取目标中心的大小为 $a^n P \times a^n R$ 的窗口标记为 J^n 。 a 表示尺度因子，取 $a = 1.02$ ， $S = 33$ ， n 的取值范围如下：

$$n \in \left\{ \left\lfloor -\frac{S-1}{2} \right\rfloor, \dots, \left\lfloor \frac{S-1}{2} \right\rfloor \right\}$$

尺度等级为 n 的训练样本 f 是 J^n 的 d 维的特征描述子。

在连续的两帧中，位置的变化往往大于尺度的变化，因此，文中先采用位置滤波器确定位置信息，在位置的基础上再使用尺度滤波器确定尺度信息。

- 算法流程

Algorithm 1 Proposed tracking approach: iteration at time step t .

Input:

Image I_t .

Previous target position \mathbf{p}_{t-1} and scale s_{t-1} .

Translation model $A_{t-1}^{\text{trans}}, B_{t-1}^{\text{trans}}$ and scale model $A_{t-1}^{\text{scale}}, B_{t-1}^{\text{scale}}$.

Output:

Estimated target position \mathbf{p}_t and scale s_t .

Updated translation model $A_t^{\text{trans}}, B_t^{\text{trans}}$ and scale model $A_t^{\text{scale}}, B_t^{\text{scale}}$.

Translation estimation:

- 1: Extract a translation sample z_{trans} from I_t at \mathbf{p}_{t-1} and s_{t-1} .
- 2: Compute the translation correlation y_{trans} using $z_{\text{trans}}, A_{t-1}^{\text{trans}}$ and B_{t-1}^{trans} in (6).
- 3: Set \mathbf{p}_t to the target position that maximizes y_{trans} .

Scale estimation:

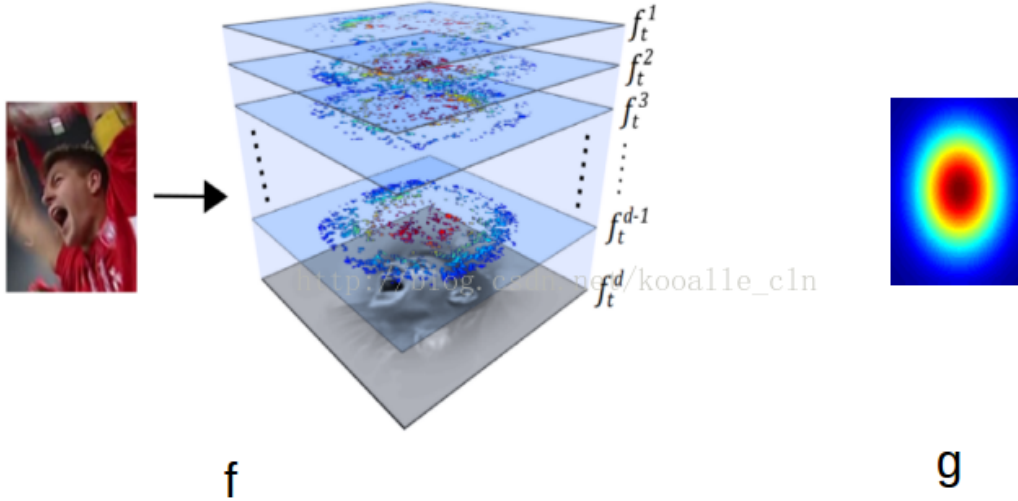
- 4: Extract a scale sample z_{scale} from I_t at \mathbf{p}_t and s_{t-1} .
- 5: Compute the scale correlation y_{scale} using $z_{\text{scale}}, A_{t-1}^{\text{scale}}$ and B_{t-1}^{scale} in (6).
- 6: Set s_t to the target scale that maximizes y_{scale} .

Model update:

- 7: Extract samples f_{trans} and f_{scale} from I_t at \mathbf{p}_t and s_t .
 - 8: Update the translation model $A_t^{\text{trans}}, B_t^{\text{trans}}$ using (5).
 - 9: Update the scale model $A_t^{\text{scale}}, B_t^{\text{scale}}$ using (5).
-

- 位置估计的训练

目标所在的图像块 P 的大小为 $M \times N$ ，提取 P 的特征（例如Fhog特征），得到大小为 $M \times N \times d$ 的特征 f ，其中特征的维度为 d 维。如下图所示（作者选取灰度图为第1维的特征，后续的特征是Fhog的前27维特征），响应 g 是根据高斯函数构造的，大小为 $M \times N$ ，中间响应值最大，向四周依次递减：



根据下述公式，需要对 f 的每一个维度的特征做二维的DFT，得到 F^l ，对 g 做二维的DFT得到 G ：
（下式中的乘法都是点乘，矩阵的对应位置相乘）

$$H^l = \frac{\bar{G} F^l}{\sum_{k=1}^d \bar{F}^k F^k + \lambda}$$

- 位置估计的检测

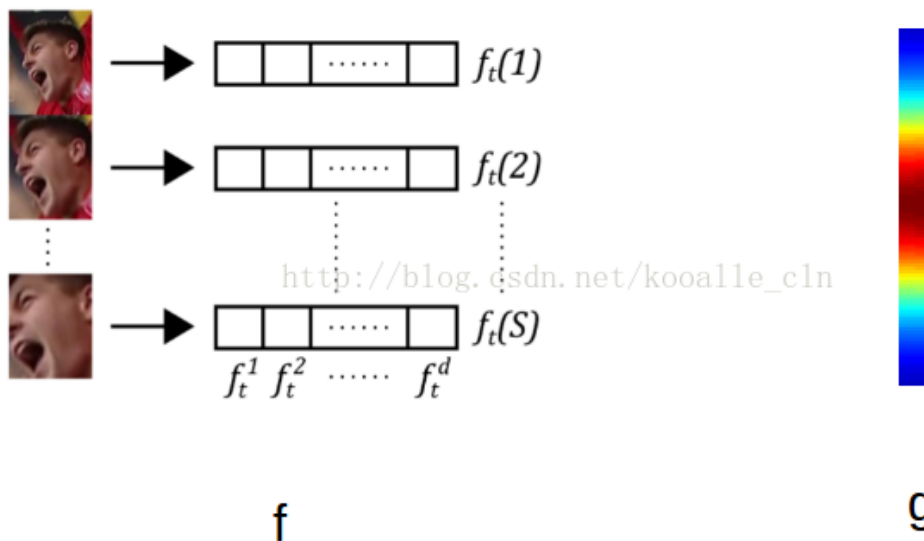
新的一帧的图片的特征为 z ，同样求取每一维度特征的DFT得到 Z^l ，按下式求出响应 y 。

$$y = \mathcal{F}^{-1} \left\{ \frac{\sum_{l=1}^d \bar{A}^l Z^l}{B + \lambda} \right\}$$

- 尺度估计的训练

目标所在的图像块 P 的大小为 $M \times N$ ，以图像块的正中间为中心，截取不同尺度的图片。这样就能够得到一系列的不同尺度的图像Patch（搜索范围为 S 个尺度，就会有 S 张图像Patch），针对每个图像Patch求其特征描述子（维度为 d 维，这里的 d 和位置估计中的维度 d 没有任何关系，只与截取图像resize的尺寸有关），每一个维度的特征 f^l 是一个 $1 \times S$ 的向量， g 是高斯函数构造的输出响应大小为 $1 \times S$ ，中间值最大，向两端依次减小。

如上一样得到滤波器模板 H



- 尺度估计的检测

新的一帧的图片以位置估计得出的位置为中心，截取 S 个不同尺度的图像Patch，分别求其特征描述子，组成新的特征为 z ，同样求取每一维度的一维的DFT得到 Z^l ，进而得到 y 的值（求取方法和位置估计类似）， y 为 $1 \times S$ 维的向量， y 中最大值的所对应的尺度为最终尺度估计的结果。

- 总结

1. Input：前一帧的目标位置 P_{t-1} 和尺度 S_{t-1} ，位置滤波器参数 A_{t-1}^{trans} 、 B_{t-1}^{trans} 和尺度滤波器参数 A_{t-1}^{scale} 、 B_{t-1}^{scale}
2. Output：估计的目标位置和尺度 P_t 、 S_t
3. 更新：利用新一帧的位置和尺度，更新滤波器参数。

4. Fast Discriminative Scale Space Tracking

4.1 Sub-Grid Interpolation of Correlation Scores

4.2 Dimensionality Reduction

为了减少傅里叶变换的计算量，更新一个目标模板：

$$u_t = (1 - \eta)u_{t-1} + \eta f_t$$

则滤波器的分子 A_t^l 可以通过 $\bar{G} \mathcal{F} \{u_t^l\}$ 得到。利用模板 u_t 构造了一个投影矩阵 P_t ($\tilde{d} \times d$)。这个矩阵定义了特征被投影到上面的低维子空间， \tilde{d} 为特征压缩后的维数。

- 通过最小化目标模板 u_t 的重构误差来获得 P_t

$$\varepsilon = \sum_n \left\| u_t(\mathbf{n}) - P_t^T P_t u_t(\mathbf{n}) \right\|^2$$

$$s.t. P_t P_t^T = I$$

- 通过对自相关矩阵进行特征值分解，得到一个解：

$$C_t = \sum_n u_t(\mathbf{n}) u_t(\mathbf{n})^T$$

- C_t 的前 \tilde{d} 个最大的特征值对应的特征向量构成了 P_t 的行向量
- 压缩维度后的训练样本：

$$\tilde{F}_t = \mathcal{F} \{P_t f_t\}$$

- 压缩目标模板：

$$\tilde{U}_t = \mathcal{F} \{P_t u_t\}$$

- FDSST更新策略：

$$\tilde{A}_t^l = \bar{G} \tilde{U}_t^l, \quad l = 1, \dots, \tilde{d}$$

$$\tilde{B}_t = (1 - \eta) \tilde{B}_{t-1} + \eta \sum_{k=1}^{\tilde{d}} \tilde{F}_t^k \tilde{F}_t^k$$

- 检测：

$$\tilde{Z}_t = \mathcal{F} \{P_{t-1} z_t\}$$

$$Y_t = \frac{\sum_{l=1}^{\tilde{d}} \overline{\tilde{A}_{t-1}^l} \tilde{Z}_t^l}{\tilde{B}_{t-1} + \lambda}$$