

---

---

# COMPUTER VISION NOTES

## OBJECT LOCALISATION AND TRACKING

---

---

MY PERSONAL NOTES ON

OBJECT LOCALISATION TECHNIQUES; COLOUR MATCHING, MEAN SHIFT TRACKING, OPTICAL FLOW,  
LUKAS KANADE

BY

0xLeo ([github.com/0xleo](https://github.com/0xleo))

APRIL 11, 2019

DRAFT X.YY

MISSING: ...

## Contents

<b>1</b>	<b>Histogram-Based Methods</b>	<b>2</b>
1.1	Histogram backprojection . . . . .	2
1.1.1	Intuition - model and search image histogram . . . . .	2
1.1.2	(optional) The rationale behind defining the ratio histogram . . . . .	2
1.1.3	A colour matching heuristic . . . . .	3
1.1.4	Histogram backprojection implementation from scratch . . . . .	4
1.1.5	Histogram backprojection implementation using OpenCV's API . . . . .	4
1.1.6	Histogram backprojection summary . . . . .	5
1.2	Mean Shift Tracking . . . . .	5
1.2.1	Mean Shift algorithm idea . . . . .	5
1.2.2	Mean Shift terminology and notation . . . . .	7
1.2.3	Mathematical analysis . . . . .	7
1.2.4	Mean Shift as a tool for segmentation . . . . .	8
1.2.5	Mean Shift as a tracker . . . . .	9
1.2.6	Implementation from scratch . . . . .	9
1.2.7	Implementation in OpenCV . . . . .	10
1.2.8	Mean Shift pros and cons . . . . .	10
1.3	Camshift (Continuously Adaptive Mean Shift) . . . . .	10
<b>2</b>	<b>Motion-based methods</b>	<b>10</b>
2.1	Optical Flow . . . . .	10
2.2	Lukas-Kanade tracking . . . . .	10
<b>A</b>	<b>Appendices</b>	<b>11</b>
A.1	HSV domain . . . . .	12
A.2	Histogram backprojection implementation from scratch – source code . . . . .	13
A.3	Histogram backprojection implementation using OpenCV - source code . . . . .	15
A.4	My mean shift implementation from scratch – source code . . . . .	17
A.5	Mean Shift implementation using OpenCV - source code . . . . .	20

# 1 Histogram-Based Methods

## 1.1 Histogram backprojection

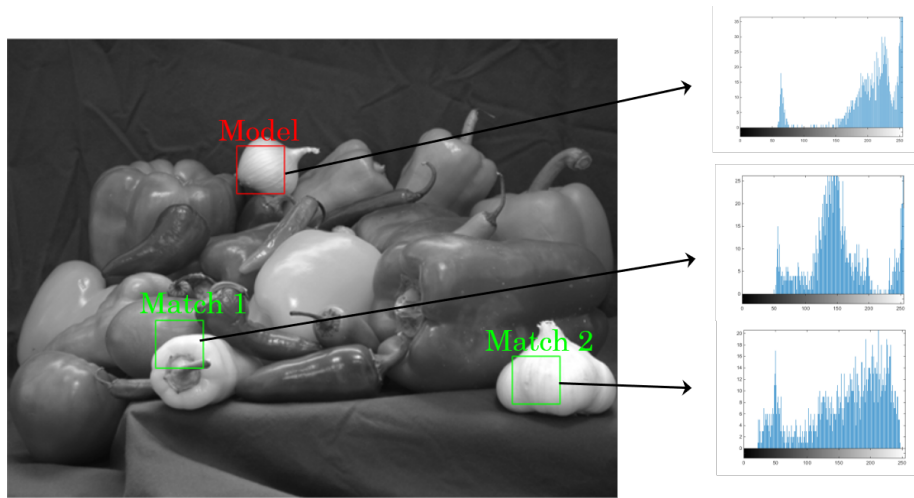
### 1.1.1 Intuition - model and search image histogram

In image processing, we are usually interested in histograms of greyscale images. However, often the colour histogram can be used to identify an image region or object. RGB histograms are practically not good enough for matching as the R, G, B components are strongly correlation with the illumination hitting the object. In practice, objects are converted from RGB to HSV (Hue, Saturation, Value) domain. `Hue` represents the colour type (blue, yellow, etc.), `saturation` represents the vibrancy (how vivid or neutral it is) and `value` represents the brightness of the colour. Hence HSV decouples the brightness from the colour description. Therefore when performing colour matching we are only interested in the H and S components, which map to a 2D histogram. More about the HSV domain in A.1.

☞ **Comment:** The HS components are often but *not always* a good choice for colour-based detection. They may fail detecting black and white objects since black and white can have any colour (H) and in this case the SV components of the HSV or even the YUV domain are a better choice. However, in this article we stick to HS.

Histogram backprojection answers the question “where in the image are the colours that belong to the object being looked for?”. We do this by defining a model image (the object we search for – a. k.a. target) and the search (the whole image where we search in), probing the model over search image and calculating their histogram similarity at each position.

Just to illustrate the idea, assume that we want to match the greyscale (instead of the 2D) histogram of the garlic in Fig. 1. A part of the top garlic has been chosen as the model. The histogram of the model is shown as well as that of two matching candidates. In this case, the histogram of “match 2” is more similar to the model’s than one “match 1” so we want somehow to register that similarity. The question attempted to be answered in the next section is “how do we measure the similarity of the histogram of the matching candidate to that of the model?”.



**Fig. 1.** Model and two matches' greyscale histograms.

### 1.1.2 (optional) The rationale behind defining the ratio histogram

We assume that

1. the model's histogram (the object we search for) is narrow and tall
2. and the scene's (whole input image) histogram is rather wide.

It has then been proven (it won't be discussed in this article how as the how is out of scope) that a good histogram similarity measure between the model  $M$  and a patch of the image we search it in  $I$  is the ratio histogram  $R$ .  $M$  and  $I$  need to be pre-calculated and can be divided element-wise, since they have the same range and bins to obtain  $R$ . The ratio histogram is a therefore function  $R : \mathbb{Z}^2 \rightarrow \mathbb{R}$  that maps a colour  $(h, s)$  to some value. If  $M(h, s) > I(h, s) \Rightarrow R(h, s) > 1$  that means the model has more pixels of colour  $(h, s)$

relative to its total number of pixels compared to  $I(h, s)$ , and vice versa. If  $R(h, s) = 1$ , then that means that the input and model images contain  $(h, s)$  at the same degree.

However, because of assumptions (1), (2),  $M(h, s) > I(h, s)$  can happen quite often and so long as  $M(h, s) > I(h, s) \Rightarrow R(h, s) > 1$ , e.g.  $R(h, s) = 2, 3, 10$ , the exact value does not give much useful information. To summarise,  $R$  de-emphasises pixels with colours that do not belong to the model and emphasises the rest.

### 1.1.3 A colour matching heuristic

From the previous section, the conclusion is that it is desirable to clip the ratio histogram to 1, as defined by Swain et al.

**DEFINITION 1.1.** For each bin  $j$ , the ratio histogram is defined as

$$R_j = \min\left(\frac{M_j}{I_j}, 1\right) \quad (1.1)$$

, where  $M, I$  are the model's and input's histograms respectively.

Note that  $j$  does not necessarily have to be a pair  $(h, s)$ , but if a histogram bin (index) is quantised it can be a rectangle in the 2D space, such as  $[20, 39] \times [50, 69]$ . As mentioned before,  $R$  associates a colour with its probability of appearing in the model and the next step is to associate each pixel with that probability.

Each pixel of the original image at  $(x, y)$  maps to a 2D HS value, by a colour function  $c : \mathbb{Z}^2 \leftarrow \mathbb{Z}^2$ , by taking  $c(x, y)$ . Sometimes need an intermediate function  $h : \mathbb{Z}^2 \rightarrow \mathbb{Z}^2$  that takes the output of  $c$  and quantises it (groups multiple colours in one bin), before it is fed to  $R$ . For example,  $h$  could convert  $[0, 1, \dots, 179] \times [0, 1, \dots, 255]$  to  $[0, 19, 39, \dots, 179] \times [0, 24, 49, \dots, 255]$ . The output of  $h$  is fed to  $R$ , which divides  $M_j$  to  $I_j$  at each bin  $j$ . To summarise this paragraph we have defined the following functions in backpropagation:

- $c : \mathbb{Z}^2 \rightarrow \mathbb{Z}^2$ : maps a pixel at  $(x, y)$  to an HS value  $(h_j, s_j)$ .
- $h : \mathbb{Z}^2 \rightarrow \mathbb{Z}^2$  maps a set of values  $(h_i, s_i, h_{i+1}, s_{i+1}, \dots, h_n, s_n)$  to another  $(h, s)$  value by having quantised the range of  $h$  and  $s$ .
- $R : \mathbb{Z}^2 \rightarrow [0, 1]$  maps an  $(h, s)$  value to a probability.

We therefore want to create a new image  $b$  where each pixel  $(x, y)$  gets assigned its output of  $R$  - the measure of how much its colour appears in the model image.

$$b(x, y) := R(h(c(x, y))) = \min\left(\frac{M(h(c(x, y)))}{I(h(c(x, y)))}, 1\right) \quad \forall x, y \quad (1.2)$$

The final step is to find compact regions where  $b$  is high. If the shape of the object (model) to detect is generic, then this can be done by convolving  $b$  with binary disk mask  $D^r$  of radius  $r$ . Define:

$$D_{x,y}^r = \begin{cases} 1 & \sqrt{x^2 + y^2} \leq r \\ 0 & \text{otherwise} \end{cases} \quad (1.3)$$

Then the probability image  $b$  can be convolved with the mask:

$$b := D^r * b \quad (1.4)$$

The arg max function returns the pixel  $(x, y)$  with the maximum value of its argument, i.e. of the  $R$  matrix and the  $*$  symbol denotes convolution. Then Histogram Backprojection algorithm can be then written

---

**Algorithm 1** Colour matching by histogram backprojection according to Swain et al

---

```
1: procedure HIST-BACKPROJ(ImM, ImI) ▷ ImM: model, ImI: search image
2:    $M \leftarrow \text{histogram}(ImM)$ 
3:    $I \leftarrow \text{histogram}(ImI)$ 
4:   for each histogram bin  $j$  do ▷ a bin is a pair  $(h, s)$ 
5:      $R_j = \min\left(\frac{M_j}{I_j}, 1\right)$  ▷ Divide element-wise
6:    $m \leftarrow \text{rows}(M)$ 
7:    $n \leftarrow \text{cols}(M)$ 
8:    $b \leftarrow \text{empty}_{m \times n}$ 
9:   for  $y$  in  $0 \dots m-1$  do
10:    for  $x$  in  $0 \dots n-1$  do
11:       $b_{x,y} \leftarrow R(h(c(x,y)))$  ▷ b matrix of colour probability
12:     $D^r \leftarrow \text{binary disk of radius } r$ 
13:     $b \leftarrow D^r * b$  ▷ Group (by convolving) high probability pixels together.
14:     $x_{obj}, y_{obj} \leftarrow \arg \max_{x,y} (b)$ 
15:  return  $x_{obj}, y_{obj}$ 
```

---

#### 1.1.4 Histogram backprojection implementation from scratch

An implementation of Alg. 1 has been written in A.2. However, instead of finding the location of the object by the argmax function, it applies Otsu's threshold on the  $R$  matrix. This automatically selects a threshold  $T$  based on the statistics of the histogram of  $R$  for which if  $R[x,y] < T$ , then the pixel at  $(x,y)$  is classified as background, else as foreground. Instructions on how to run the implementation code are in A.2 and an output is shown below.



**Fig. 2.** Input image with a ROI of the objects to detect selected.



**Fig. 3.** Detected objects on the original image.

#### 1.1.5 Histogram backprojection implementation using OpenCV's API

OpenCV implements the technique using the `cv2.calcBackProject(image, channels, histogram_array, channel_ranges, [scale = 1])` method (in Python). Its invocation looks like:

`cv2.calcBackProject(search_image, channels, model_histogram, channel_ranges, [scale = 1])`

- `search_image`: the input image, e.g. in HSV.
- `channels`: which channels of the original image and the model to select in order to draw its histogram, e.g. `channels = [0,1]` -> H, S.
- `model_histogram`: histogram of the model (ROI), needs to be pre-calculated.
- `channel_ranges`: set it to `[0,180,0,256]` to select the full range of H, S components.

The code listing in A.3 works similarly with the one in A.2, expecting two clicks from the user to define a bounding box around a sample of the object to detect. It also performs similarly on the same images, showing some black spots on roughly the same positions.

### 1.1.6 Histogram backprojection summary

- ✓ Fast – can easily be used in real time.
- ✓ Relatively immune to noise and illumination changes.
- ✓ Simple to implement.
- ✗ Not effective against non-compact objects.
- ✗ Does not use any knowledge about the shape or position of the detected object – only its colour.

## 1.2 Mean Shift Tracking

### 1.2.1 Mean Shift algorithm idea

Mean shift is a non-parametric feature-space analysis technique for locating the maxima of a density function, a so-called mode-seeking algorithm. In image processing, it's used for tracking.

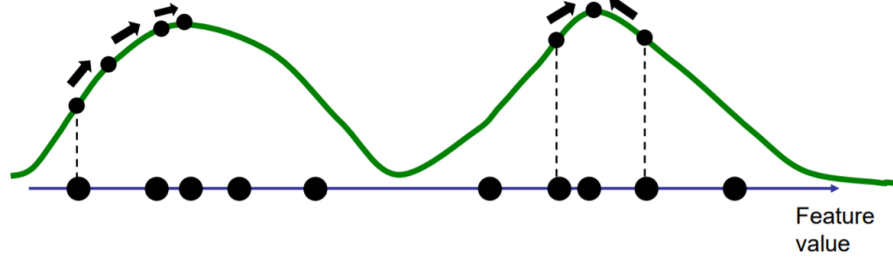
Given some points  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  with  $d$  dimensions (“features”), for each data point, mean shift defines a window around it and computes the mean of data point. Then it shifts the centre of window to the mean and repeats the algorithm till the window stops moving. In image processing, feature space is often the colour space. Table 1 illustrates the iterations until the algorithm converges.

Mean shift is a nonparametric iterative algorithm. It considers each point sampled from a probability distribution, i.e. each point is most likely found at its actual measured position, but it can also be in a neighbourhood around it.

Table 1: Mean shift update steps shown on a very high level – in this case the new centre is simply the centroid, until converge (last row).

Current centre (blue)	Mean shift vector	New centre (yellow)

- ❓ How do we find the peak towards which the circle should move?
- ✅ The circle should move towards the densest point of the distribution of all points within the ROI. The peak is found by superimposing all the individual probability distributions around each point and finding the  $N$  highest maxima of the result ( $N$  is the number of classes we want to have).
- ❓ How do we convert a set of discrete input points to a continuous density function so that we can find the maxima (Fig. 4)?



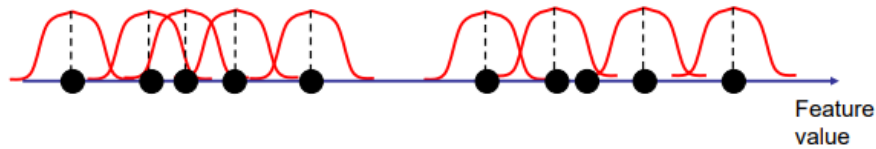
**Fig. 4.** The green curve is what we roughly want to generate from the 1D input points. The arrows simply show the path gradient ascent would follow to locate the maxima.

- ✅ Let us define a kernel function.

**DEFINITION 1.2.** A kernel is a real-valued function of the points  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  that satisfies the following properties:

1.  $K$  is maximum at  $\mathbf{0}$ , non increasing, and decays away from the maximum.
2.  $K$  is radially symmetric.
3.  $K(\mathbf{x}) \geq 0$ .
4.  $\int_{\mathbb{R}^d} K(\mathbf{x}) d\mathbf{x} = 1$

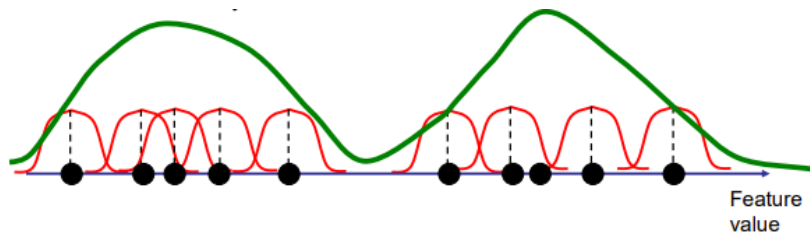
Then for each input point, its kernel function should look roughly as follows.



**Fig. 5.** The kernel functions  $K(\mathbf{x} - \mathbf{x}_i)$ , where  $\mathbf{x}_i$  are the input points.

If we allocate each point  $\mathbf{x}_i$  its own kernel  $K(\mathbf{x} - \mathbf{x}_i)$ , then by summing all  $N$  of them and dividing by a constant  $C$  to normalise the result we can get a probability density function (PDF):

$$f(\mathbf{x}) = \frac{1}{C} \sum_{i=1}^N K(\mathbf{x} - \mathbf{x}_i) \quad (1.5)$$



**Fig. 6.** The sum of individual kernel functions functions  $K(\mathbf{x} - \mathbf{x}_i)$ .

$f(\mathbf{x})$  approximates the probability that feature  $\mathbf{x}$  is observed given the data points. The maxima of  $f$  (the “modes” of the pdf) correspond to the clusters in the data. As shown in Fig. 4, a way to reach the peak of the PDF  $f$  is by incrementing the mean by  $\nabla f(\mathbf{x})$ . A very rough algorithm for Mean Shift would therefore be.

---

**Algorithm 2** Mean shift on a very high level

---

```

1: procedure MEANSHIFTIDEA( $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ )
2:   for  $i = 1, \dots, N$  do
3:      $\mathbf{x} \leftarrow \mathbf{x}_i$ 
4:     while no convergence do ▷ gradient ascent for each point individually
5:        $\mathbf{x} \leftarrow \mathbf{x}_i + \nabla f(\mathbf{x}) = \mathbf{x}_i + \frac{1}{C} \sum_i \nabla K(\mathbf{x} - \mathbf{x}_i)$ 
6:   return  $\mathbf{x}$ 

```

---

### 1.2.2 Mean Shift terminology and notation

Before the maths is presented, this is notation that will be used.

- $d$  – the dimension of the input column vector, the entries of this vector are also called “features”.
- $\mathbf{x}_i$  – the data points.
- “Kernel”  $K(\mathbf{x})$  – the function that assigns weight to every point of interest. For example, it can be Gaussian, flat, etc.
- “Bandwidth”  $h$  – the radius of the region of interest (ROI).
- “PDF”  $f(\mathbf{x})$  – probability density function.

### 1.2.3 Mathematical analysis

A basic requirement for the kernel function, as stated in Def. 1.2 is radial symmetry, i.e.

$$K(\mathbf{x}) = c_d k(\|\mathbf{x}\|^2) \quad (1.6)$$

, where  $c_d$  acts as the normalisation constant and is the volume of the  $d$ -dimensional sphere, such that  $K(\mathbf{x})$  integrates to 1. Some functions that can serve as the basis  $k$  for the kernel are

$$\text{Epanechnikov} \quad k_E(\mathbf{x}) = \begin{cases} 1 - \|\mathbf{x}\|^2 & \|\mathbf{x}\| \leq 1 \\ 0 & \text{otherwise} \end{cases} \quad (1.7)$$

$$\text{Uniform} \quad k_U(\mathbf{x}) = \begin{cases} 1 & \|\mathbf{x}\| \leq 1 \\ 0 & \text{otherwise} \end{cases} \quad (1.8)$$

$$\text{Gaussian} \quad k_N(\mathbf{x}) = \exp\left(-\frac{1}{2} \|\mathbf{x}\|^2\right) \quad (1.9)$$

It can be proven that the pdf function that approximates kernel density given inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , where  $\mathbf{x}_i \in \mathbb{R}^d$  is

$$\hat{f}(\mathbf{x}) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right) \quad (1.10)$$

Assuming that the kernel  $K(\mathbf{x})$  is differentiable, the gradient of the kernel density is

$$\hat{\nabla} f(\mathbf{x}) := \nabla \hat{f}(\mathbf{x}) = \frac{1}{h^d} \sum_{i=1}^n \nabla K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right) \quad (1.11)$$

Computing the gradient, the derivative is

$$\nabla \hat{f}(\mathbf{x}) = \frac{2c_d}{nh^{d+2}} \sum_{i=1}^n (\mathbf{x} - \mathbf{x}_i) k'\left(\left\|\frac{\mathbf{x} - \mathbf{x}_i}{h}\right\|^2\right) \quad (1.12)$$

$$= \frac{2c_d}{nh^{d+2}} \left( \sum_{i=1}^n g_i \right) \underbrace{\left( \frac{\sum_{i=1}^n \mathbf{x}_i g_i}{\sum_{i=1}^n g_i} - \mathbf{x} \right)}_{\text{mean shift vector}}, \quad (1.13)$$

$$g(r) = k'(r), \quad r := \|\mathbf{x}\|, \quad g_i = g(\|\mathbf{x} - \mathbf{x}_i\|/h)^2 \quad (1.14)$$



**DEFINITION 1.3 (mean shift vector).** Referring to Eq. (1.14),  $\sum_{i=1}^n g_i$  is yet another kernel estimation and the second term  $M(\mathbf{x}) = \frac{\sum_{i=1}^n \mathbf{x}_i g_i}{\sum_{i=1}^n g_i} - \mathbf{x}$  is the mean shift vector.

It always points toward the direction of the maximum increase in the density therefore we want to shift the current estimation  $\mathbf{x}$  by it in each iteration. The mean shift vector must be  $\mathbf{0}$  at optimum, i.e. the algorithm stops at  $\mathbf{x} = \frac{\sum_{i=1}^n \mathbf{x}_i g_i}{\sum_{i=1}^n g_i}$ . This is equivalent to draft Alg. 2 converging. Given this knowledge about the gradient, the latter algorithm is rewritten as follows.

---

**Algorithm 3** Mean shift algorithm

---

```

1: procedure MEANSHIFT( $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ )
2:   for  $i = 1, \dots, n$  do
3:      $\mathbf{x} \leftarrow \mathbf{x}_i$ 
4:     while no convergence do ▷ gradient ascent for each point individually
5:        $\mathbf{x} \leftarrow \mathbf{x} + M(\mathbf{x}) = \frac{\sum_{i=1}^n \mathbf{x}_i g\left(\left\|\frac{\mathbf{x}-\mathbf{x}_i}{h}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{\mathbf{x}-\mathbf{x}_i}{h}\right\|^2\right)}, \quad g(\|\mathbf{x}\|) = k'(\|\mathbf{x}\|)$ 
6:   return  $\mathbf{x}$ 

```

---

guaranteed to converge. The update step is not as complicated as it seems. For example for the kernel

$$k(\|\mathbf{x}\|) = \begin{cases} 1 - \|\mathbf{x}\| & \|\mathbf{x}\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

, its derivative w.r.t. the distance  $\|\mathbf{x}\|$  is

$$g(\|\mathbf{x}\|) = \begin{cases} -1 & \|\mathbf{x}\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Comparing the argument of  $g$  with 1 gives us the data points of interest therefore the “mean” part of  $M(\mathbf{x})$  is:

$$\frac{\sum_{i=1}^n \mathbf{x}_i g\left(\left\|\frac{\mathbf{x}-\mathbf{x}_i}{h}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{\mathbf{x}-\mathbf{x}_i}{h}\right\|^2\right)} = \frac{-\sum_{\|\mathbf{x}-\mathbf{x}_i\| < h} \mathbf{x}_i}{-\sum_{\|\mathbf{x}-\mathbf{x}_i\| < h} 1} = \frac{\sum_{\|\mathbf{x}-\mathbf{x}_i\| < h} \mathbf{x}_i}{n_h} \quad (1.15)$$

$n_h$  is simply the number of inside the kernel, for which  $\|\mathbf{x} - \mathbf{x}_i\| < h$  and  $\sum_{\|\mathbf{x}-\mathbf{x}_i\| < h} \mathbf{x}_i$  is just the average of the data points within a radius  $h$  of  $\mathbf{x}$ ! Regarding the convergence of the algorithm the following can be proved.

**THEOREM 1.1.** *If the kernel function  $k(\mathbf{x})$  is convex and monotonically decreasing then the update of  $\mathbf{x}$  converges and the pdf  $\hat{f}(\mathbf{x})$  increases.*

For the Epanechnikov kernel, convergence is reached in finite number of steps. Finally, mean shift runs in  $\mathcal{O}(n^2 T)$ , where  $n$  is the number of input points and  $T$  the number of iterations.

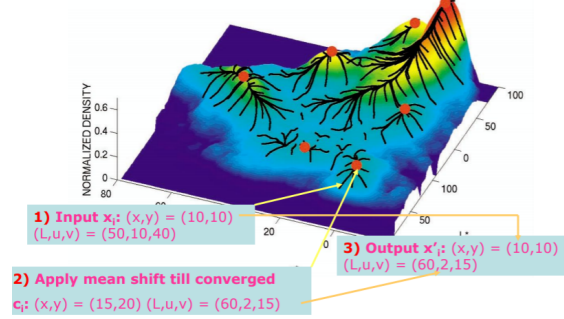
#### 1.2.4 Mean Shift as a tool for segmentation

When we want to segment an image, it is usually converted from RGB to another colour space such as HSV or LUV. In this case, suppose the image is in LUV.

Then the feature space is  $(l, u, v, x, y)$ . To perform segmentation, we apply *two different* mean shifts in the 5-dimensional space as we want to segment pixels based on their location and colour. For each pixel  $(x_i, y_i)$  of intensity color  $(l_i, u_i, v_i)$ , find the corresponding mode  $c_{col}$ . All of the pixels  $(x_i, y_i)$  corresponding to the same mode  $c_{col}$  are grouped into a single region. At the same time, mean shift is performed in the 2D space  $xy$  and all of the corresponding pixels  $(x_i, y_i)$  are grouped into a single mode  $c_{pos}$ . The kernel in this case is

the product of the position kernel and the colour kernel,

$$K_{pos,col}(\mathbf{x}) = \frac{c}{h_{col}^3 h_{pos}^2} k\left(\frac{\|\mathbf{x}_{pos}\|^2}{h_{pos}^2}\right) k\left(\frac{\|\mathbf{x}_{col}\|^2}{h_{col}^2}\right) \quad (1.16)$$



**Fig. 7.** Mean shift in the LUV space.

### 1.2.5 Mean Shift as a tracker

The idea behind using it as a tracker is the following. Given a frame that contains the object of interest (OOI), we want to successively move the window of radius  $h$  towards the “centre” of the object. How is the centre defined?

After selecting a sample of the OOI, backprojection from the previous section helps create a greyscale. salience (likelihood) frame (see  $b$  matrix in backprojection section). In this frame, the whiter the pixel, the more likely it is to belong to the OOI.

Then, mean shift can be performed in 3D -  $(x, y, i)$ , where  $i$  is the salience intensity to chase the object and output its  $(x, y)$ .

### 1.2.6 Implementation from scratch

Below are the main features of my own adaptation of mean shift.

- Mean shift is performed in the  $xy$  space so the “mode” (converge) centre is 2D.
- If it performed in  $xy$  space, which points are considered? The candidate points are the result of backprojecting a sample of the OOI to the frame, create a likelihood frame  $b$ , as described in Section 1.1. Then they are automatically thresholded, which yields a BW frame. For that frame, the white points are likely to belong to the OOI and the black are irrelevant. So we want to generate a pdf from those white points and find its centre of density.
- The chosen mean shift kernel is  $k(r) = 1 - r$ ,  $r \leq 1$ ,  $r = \|x\|$ . The radius of interest can be chosen by the user but defaulted to some small value anyway. Therefore the update step of the m.s. vector  $x$  is simple – it is simply the centroid of all points up to distance  $h$  around  $x$  as derived in 1.15.

My implementation code is listed in A.2. The program is standalone and expect a video path and a radius from the user. Below are some representative outputs for a football sequence. For that particular scene, the player was successfully tracked for the majority of the frames.



**Fig. 8.** Initial step; grabbing a sample of the player to track (green box).



**Fig. 9.** An intermediate tracking step (blue circle).



**Fig. 10.** Player keeps being tracked right after being blocked by another player.

The algorithm is robust given that the sample represents reasonably well the OOI. It is not robust when the scale of the OOI changes, e.g. in car scenes when a car to be tracked car is initially in the background and in the end approaches the camera. Then mean shift variations that are able to adapt the radius  $h$  need to be considered and there's a good amount of literature on that.

### 1.2.7 Implementation in OpenCV

Mean shift in OpenCV consists of the following main stages:

1. Set up a target; i.e. grab a sample of the OOI. Once again, the OOI shouldn't be in RGB, but in a colour space that decouples hue and illumination.
2. Provide an initial window – essentially initialise the mean shift tracking vector  $x$  and provide a radius  $h$  to define the area around  $x$  where points for the update of  $x$  will be considered.
3. Perform backprojection of the target on the current frame, generating a greyscale image. Perform mean shift on the greyscale backprojected image.

After grabbing the ROI (in HSV), the OpenCV approach also filters only certain H (0 to 180 – maximum is 180), S (60 to 255 – maximum is 255), V (32 to 255 – maximum is 255) bands in order to remove potential noise:

```
mask = cv2.inRange(hsv_roi, np.array((0., 60.,32.)), np.array((180.,255.,255.)))
```

Mean shift is performed, for instance, as:

```
cv2.meanShift(probImage, track_window, criteria) -> retval, track_window
```

- probImage: greyscale backprojected frame.
- track\_window: the updated tracking window as  $x$ ,  $y$ ,  $w$ ,  $h$  (essentially the updated m.s. vector).
- term\_crit: termination criteria, e.g. to set them to 10 iterations or window displacement no more than 1 pixel do ( cv2.TERM\_CRITERIA\_EPS | cv2.TERM\_CRITERIA\_COUNT, 10, 1).
- retval: A boolean whether the algorithm was successful.
- track\_window: The updated tracking window in the same format as the old one.

The code provided by OpenCV's documentation has been lazily re-written and made interactive as listed in A.5. The user can define the bounding box around the object of interest and then press q to finalise the selection and start the algorithm.

### 1.2.8 Mean Shift pros and cons

- |  |  |
|--|--|
| ✓ Does not assume spherical clusters.    | ✗ Output depends on window size.   |
| ✓ Just a single parameter (window size). | ✗ Computationally expensive (however, not impossible to use in real time). |
| ✓ Finds variable number of modes.        | ✗ Does not scale well with dimension of feature space.                     |
| ✓ Robust to outliers.                    |  |

## 1.3 Camshift (Continuously Adaptive Mean Shift)

## 2 Motion-based methods

### 2.1 Optical Flow

### 2.2 Lukas-Kanade tracking

**A   Appendices**

## A.1 HSV domain

## A.2 Histogram backprojection implementation from scratch – source code

**Listing 1:** Histogram backprojction from scratch (src/hist\_backproj/backproj.py).

```
1 import cv2, numpy as np
2 import sys
3 import pdb
4
5 g_clicks_xy = []
6
7 def qimshow(im, delay = 10, wname = 'display'):
8     cv2.imshow(wname, im)
9     cv2.waitKey(delay * 1000)
10    cv2.destroyAllWindows()
11
12
13 def on_click(event, x, y, flags, param):
14     """
15     Mouse callback function - write to global list of clicks
16     """
17     global g_clicks_xy
18     if event == cv2.EVENT_LBUTTONDOWN:
19         g_clicks_xy.append((x, y))
20
21 """
22 @im: The input image as read (in BGR)
23 """
24 def get_model(im):
25     """
26     Process user input (clicks) and Extract the target image
27     (model) in hsv.
28     """
29     global g_clicks_xy
30     hsv = cv2.cvtColor(im, cv2.COLOR_BGR2HSV)
31     cv2.namedWindow('input')
32     cv2.setMouseCallback('input', on_click)
33     while True:
34         cv2.imshow('input', im)
35         k = cv2.waitKey()
36         if k == ord('q'):
37             break
38     cv2.destroyAllWindows()
39     # for clicks, index 0 = x, index 1 = y
40     g_clicks_xy = sorted(g_clicks_xy,
41                          key = lambda x: x[0]**2 + x[1]**2,
42                          reverse = True) [-2:]
43     click_br, click_tl = g_clicks_xy[0], g_clicks_xy[1]
44     w = click_br[0] - click_tl[0]
45     h = click_br[1] - click_tl[1]
46     hsvt = hsv[click_tl[1]: click_tl[1] + h,
47              click_tl[0]: click_tl[0] + w]
48     rect = cv2.rectangle(im.copy(), g_clicks_xy[0], g_clicks_xy[1],
49                          color = (0,255,0))
50     qimshow(rect, wname = 'selection', delay = 3)
51     return hsvt
52
53 """
54 @hsv: The whole input (search) image in hsv
55 @hsvt: The target (model), i.e. the ROI, in hsv
56 @<return>: The ration histogram of hsvt by hsv, clipped from 0 to 1
57 """
58
59 def ratio_histogram(hsv, hsvt):
60     # see doc: HS histograms, [0, 180] as in OpenCV 0 <= H <= 179
61     M = cv2.calcHist([hsv], [0, 1], None, [180, 256], [0, 180, 0, 256])
```

```

62     I = cv2.calcHist([hsvt], [0, 1], None, [180, 256], [0, 180, 0, 256])
63     R = np.divide(np.array(M, np.float),
64                   np.array(I, np.float),
65                   out = np.zeros_like(I),
66                   where = I != 0)
67     R[R > 1.0] = 1.0
68     return R
69
70
71 """
72 @hsv: the whole input image in HSV
73 @R: the ratio histogram as returned by the ratio_histogram function
74 @r: the radius of the disk backprojection convolves with
75 """
76 def backproject(hsv, R, rad = 15):
77     """
78     Generate a 2D binary image where ones are probably the object(s)
79     of interest and zeros the background
80     """
81     b = np.zeros((hsv.shape[0], hsv.shape[1]), np.uint8)
82     for r in range(hsv.shape[0]):
83         for c in range(hsv.shape[1]):
84             b[r,c] = R[hsv[r,c][0], hsv[r,c][1]]
85     b = np.uint8(b)
86     disk = cv2.getStructuringElement(cv2.MORPH_ELLIPSE, (rad, rad))
87     # convolve with disk
88     b = cv2.filter2D(b, -1, disk, b)
89     b = np.uint8(b)
90     b = np.array(cv2.normalize(b, b, 0, 255, cv2.NORM_MINMAX), np.uint8)
91     _, thresh = cv2.threshold(np.uint8(b), 0, 255,
92                               cv2.THRESH_BINARY | cv2.THRESH_OTSU )
93     return thresh
94
95
96 def usage():
97     str_usage = 'Run with $ python <program_name> <input_image>\n'\
98     'When your input image shows up, click 2 times to define\n'\
99     'a bounding box around a sample of your object of interest.\n'\
100     'Then press "q" to finish the selection.'
101     print str_usage
102     sys.exit(0)
103
104
105 def main():
106     assert len(sys.argv) > 1,\
107         "Run the program with -h or --help to print usage"
108     if sys.argv[1] == '-h' or sys.argv[1] == '--help':
109         usage()
110     else:
111         im = cv2.imread(sys.argv[1])
112         hsv = cv2.cvtColor(im, cv2.COLOR_BGR2HSV)
113         hsvt = get_model(im)
114         R = ratio_histogram(hsv, hsvt)
115         thresh = backproject(hsv, R)
116         # final processing - AND the 2D binary threshold image with the
117         # original RGB input image to show the result
118         thresh = cv2.merge((thresh, thresh, thresh))
119         res = cv2.bitwise_and(im, thresh)
120         qimshow(res)
121
122
123 if __name__ == '__main__':
124     main()

```

### A.3 Histogram backprojection implementation using OpenCV - source code

**Listing 2:** Histogram backprojction using OpenCV's calcBackProject (src/hist\_backproj/backproj\_cv.py).

```
1 import cv2
2 import numpy as np
3 import sys
4
5 g_clicks_xy = []
6
7 def qimshow(im, delay = 10, wname = 'display'):
8     cv2.imshow(wname, im)
9     cv2.waitKey(delay * 1000)
10    cv2.destroyAllWindows()
11
12
13 def on_click(event, x, y, flags, param):
14     """
15     Mouse callback function - write to global list of clicks
16     """
17     global g_clicks_xy
18     if event == cv2.EVENT_LBUTTONDOWN:
19         g_clicks_xy.append((x, y))
20
21
22 """
23 @im: The input image as read (in BGR)
24 """
25 def get_model(im):
26     """
27     Process user input (clicks) and Extract the target image
28     (model) in hsv.
29     """
30     global g_clicks_xy
31     hsv = cv2.cvtColor(im, cv2.COLOR_BGR2HSV)
32     cv2.namedWindow('input')
33     cv2.setMouseCallback('input', on_click)
34     while True:
35         cv2.imshow('input', im)
36         k = cv2.waitKey()
37         if k == ord('q'):
38             break
39     cv2.destroyAllWindows()
40     # for clicks, index 0 = x, index 1 = y
41     g_clicks_xy = sorted(g_clicks_xy,
42                          key = lambda x: x[0]**2 + x[1]**2,
43                          reverse = True) [-2:]
44     click_br, click_tl = g_clicks_xy[0], g_clicks_xy[1]
45     w = click_br[0] - click_tl[0]
46     h = click_br[1] - click_tl[1]
47     hsvt = hsv[click_tl[1]: click_tl[1] + h,
48              click_tl[0]: click_tl[0] + w]
49     rect = cv2.rectangle(im.copy(), g_clicks_xy[0], g_clicks_xy[1],
50                          color = (0,255,0))
51     qimshow(rect, wname = 'selection', delay = 3)
52     return hsvt
53
54
55 def backproject(hsv, hsvt, rad = 9):
56     # Calculating object histogram
57     # In OpenCV, hue lies within [0,179]!
58     hsvt_hist = cv2.calcHist([hsvt], # image
59                              [0, 1], # channel selection
60                              None, # mask
61                              [90, 128], # no of bins
```



```

62         [0, 180, 0, 256]          # channel ranges
63     )
64     # normalize histogram to [0,255] and apply backprojection
65     # to get R (ratio histogram) matrix -> R between 0 and 1
66     cv2.normalize(hsvt_hist, hsvt_hist, 0, 255, cv2.NORM_MINMAX)
67     R = cv2.calcBackProject([hsv], # image
68         [0,1],                    # channel selection
69         hsvt_hist,                # histogram array
70         [0,180,0,256],           # channel ranges
71         scale = 1)
72     # convolve with circular disc
73     disc = cv2.getStructuringElement(cv2.MORPH_ELLIPSE, (rad, rad))
74     cv2.filter2D(R, -1, disc, R)
75     # Otsu's threshold
76     _, R_thresh = cv2.threshold(R, 0, 255, cv2.THRESH_BINARY | cv2.THRESH_OTSU
77 )
78     # Make it 3D to AND it with the search image
79     R_thresh = cv2.merge((R_thresh, R_thresh, R_thresh))
80     return R_thresh
81
82 def usage():
83     str_usage = 'Run with $ python <program_name> <input_image>\n'\
84     'When your input image shows up, click 2 times to define\n'\
85     'a bounding box around a sample of your object of interest.\n'\
86     'Then press "q" to finish the selection.'
87     print str_usage
88     sys.exit(0)
89
90
91 def main():
92     assert len(sys.argv) > 1,\
93         "Run the program with -h or --help to print usage"
94     if sys.argv[1] == '-h' or sys.argv[1] == '--help':
95         usage()
96     else:
97         im = cv2.imread(sys.argv[1])
98         cv2.namedWindow('select area then press q')
99         cv2.setMouseCallback('select area then press q', on_click)
100         while True:
101             cv2.imshow('select area then press q', im)
102             k = cv2.waitKey()
103             if k == ord('q'):
104                 break
105         cv2.destroyAllWindows()
106         hsv = cv2.cvtColor(im, cv2.COLOR_BGR2HSV)
107         click_br, click_tl = g_clicks_xy[0], g_clicks_xy[1]
108         w = click_br[0] - click_tl[0]
109         h = click_br[1] - click_tl[1]
110         rbgt = im[click_tl[1]: click_tl[1] + h,
111             click_tl[0]: click_tl[0] + w]
112         hsvt = cv2.cvtColor(rbgt, cv2.COLOR_BGR2HSV)
113         R_thresh = backproject(hsv, hsvt)
114         res = cv2.bitwise_and(im, R_thresh)
115         qimshow(res)
116
117
118 if __name__ == '__main__':
119     main()

```

## A.4 My mean shift implementation from scratch – source code

**Listing 3:** Mean shift using OpenCV's meanShift (src/mean\_shift/my\_mean\_shift.py).

```
1 from __future__ import print_function
2 import cv2
3 import numpy as np
4 import sys
5
6
7 g_clicks_xy = []
8
9 ##### UI #####
10 def qimshow(im, delay = 10, wname = 'display'):
11     cv2.imshow(wname, im)
12     cv2.waitKey(delay * 1000)
13     cv2.destroyAllWindows()
14
15
16 def on_click(event, x, y, flags, param):
17     """
18     Mouse callback function - write to global list of clicks
19     """
20     global g_clicks_xy
21     if event == cv2.EVENT_LBUTTONDOWN:
22         g_clicks_xy.append((x, y))
23
24 ##### Maths #####
25 def is_in_circle(centre, r, point):
26     return (centre[0] - point[0])**2 + (centre[1] - point[1])**2 <= r**2
27
28 def centroid(pts2D):
29     x, y = zip(*pts2D)
30     N = len(x)
31     return int(sum(x)/N), int(sum(y)/N)
32
33 def dist(x, y):
34     x, y = np.asarray(x), np.asarray(y)
35     return np.linalg.norm(x - y)
36
37 ##### Mean Shift process #####
38 def backproject(hsv, hsvt, rad = 9):
39     # Calculating object histogram
40     # In OpenCV, hue lies within [0,179]!
41     hsvt_hist = cv2.calcHist([hsvt], # image
42                             [0, 1], # channel selection
43                             None, # mask
44                             [90, 128], # no of bins
45                             [0, 180, 0, 256] # channel ranges
46                             )
47     # normalize histogram to [0,255] and apply backprojection
48     # to get R (ratio histogram) matrix -> R between 0 and 1
49     cv2.normalize(hsvt_hist, hsvt_hist, 0, 255, cv2.NORM_MINMAX)
50     R = cv2.calcBackProject([hsv], # image
51                             [0,1], # channel selection
52                             hsvt_hist, # histogram array
53                             [0,180,0,256], # channel ranges
54                             scale = 1)
55     # convolve with circular disc
56     if rad:
57         disc = cv2.getStructuringElement(cv2.MORPH_ELLIPSE, (rad, rad))
58         cv2.filter2D(R, -1, disc, R)
59         # Otsu's threshold
60         _, R_thresh = cv2.threshold(R, 0, 255, cv2.THRESH_BINARY | cv2.THRESH_OTSU)
61     # Make it 3D to AND it with the search image
```

```

62     #R_thresh = cv2.merge((R_thresh, R_thresh, R_thresh))
63     return R_thresh
64
65     """
66     @bw: thresholded backprojected frame (M x N)
67     @x: current mean shift vector estimation (tuple)
68     @h: mean shift ROI radius (int)
69     @<return>: new mean shift vector estimation (tuple)
70     """
71     def my_mean_shift(bw, x, h, conv_thresh = 2.0):
72         """
73         This mean shift works as follows:
74         @bw is black and white image obtained as the result of thesholding
75         the backprojected frame with the model, model being a small sample
76         of the object of interest extracted by the user.
77         The domain is simply the xy, and for the update of m.s. vector x only
78         white pixels up to h pixels around the current x are considered. White
79         pixels are most likely to be part of the ROI. The updated x is returned
80         """
81         x_old = (np.inf, np.inf)
82         while dist(x_old, x) > conv_thresh:
83             started = True
84             x_old = x
85             points_in = []
86             for r in range(bw.shape[0]):
87                 for c in range(bw.shape[1]):
88                     if bw[r, c] and is_in_circle(x, h, (c, r)):
89                         points_in.append((c, r))
90                         x = centroid(points_in)
91             print (x)
92             #qimshow(cv2.circle(bw.copy(), x, 6, 126, 4), delay = 0)
93             #qimshow(bw)
94             return x
95
96     """
97     @im: The input image as read (in BGR)
98     """
99     def get_model(im):
100         """
101         Process user input (clicks) and Extract the target image
102         (model) in hsv.
103         """
104         global g_clicks_xy
105         hsv = cv2.cvtColor(im, cv2.COLOR_BGR2HSV)
106         cv2.namedWindow('2 clicks, then press q')
107         cv2.setMouseCallback('2 clicks, then press q', on_click)
108         while True:
109             cv2.imshow('2 clicks, then press q', im)
110             k = cv2.waitKey()
111             if k == ord('q'):
112                 break
113         cv2.destroyAllWindows()
114         # for clicks, index 0 = x, index 1 = y
115         g_clicks_xy = sorted(g_clicks_xy,
116                             key = lambda x: x[0]**2 + x[1]**2,
117                             reverse = True) [-2:]
118         click_br, click_tl = g_clicks_xy[0], g_clicks_xy[1]
119         w = click_br[0] - click_tl[0]
120         h = click_br[1] - click_tl[1]
121         hsvt = hsv[click_tl[1]: click_tl[1] + h,
122                  click_tl[0]: click_tl[0] + w]
123         rect = cv2.rectangle(im.copy(), g_clicks_xy[0], g_clicks_xy[1],
124                             color = (0, 255, 0))

```

```

125     qimshow(rect, wname = 'selection', delay = 3)
126     return hsvt
127
128
129 def usage():
130     str_usage = \
131     "Instructions:\n\
132     Run the program with\n\
133     $ python <prog_name> <video> <roi_radius>\n\
134     When the \"input\" window appears, click 2 times to create a \n\
135     bounding box around a sample of the object to be tracked.\n\
136     Then press \"q\" to continue.\n"
137     print(usage)
138
139
140 ##### Driver #####
141 def main():
142     """
143     Instructions:
144     Run the program with
145     $ python <prog_name> <video> <roi_radius>"
146     When the "input" window appears, click 2 times to create a
147     bounding box around a sample of the object to be tracked.
148     Then press "q" to continue.
149     """
150     assert len(sys.argv) > 1,\
151         "\nUsage: $ python <prog_name> <video> <roi_radius>.\n\
152         Run with $ python -h to print detailed instructions"
153     if sys.argv[1] in ['-h', '--help']:
154         usage()
155
156     vid_path = sys.argv[1]
157     h = 20 if len(sys.argv) == 2 else int(sys.argv[2])
158     cap = cv2.VideoCapture('soccer2.mp4')
159     _, frame1 = cap.read()
160     # Get model (target) image
161     hsv = cv2.cvtColor(frame1, cv2.COLOR_BGR2HSV)
162     hsvt = get_model(frame1)
163
164     bw = backproject(hsv, hsvt, rad = 0)
165
166     # Initialise mean shift (x) vector from user input
167     click_br, click_tl = g_clicks_xy[0], g_clicks_xy[1]
168     w = click_br[0] - click_tl[0]
169     h = click_br[1] - click_tl[1]
170     x = (click_tl[0] + w/2, click_tl[1] + h/2)
171     valid = True
172     # frame by frame processing
173     while(valid):
174         valid, im = cap.read()
175         hsv = cv2.cvtColor(im, cv2.COLOR_BGR2HSV)
176         bw = backproject(hsv, hsvt, rad = 0)
177         x = my_mean_shift(bw, x, h)
178         cv2.imshow('detection', cv2.circle(im.copy(), x, 1, (0, 255, 0),
179         6))
180         if cv2.waitKey(1) & 0xFF == ord('q'):
181             break
182     # Release the capture
183     cap.release()
184     cv2.destroyAllWindows()
185
186
187 main()

```

## A.5 Mean Shift implementation using OpenCV - source code

**Listing 4:** Mean shift using OpenCV's meanShift (src/mean\_shift/mean\_shift\_cv.py).

```
1 import numpy as np
2 import cv2
3 import sys
4
5 g_clicks_xy = []
6
7 ##### UI #####
8 def qimshow(im, delay = 10, wname = 'display'):
9     cv2.imshow(wname, im)
10    cv2.waitKey(delay * 1000)
11    cv2.destroyAllWindows()
12
13
14 def on_click(event, x, y, flags, param):
15     """
16     Mouse callback function - write to global list of clicks
17     """
18     global g_clicks_xy
19     if event == cv2.EVENT_LBUTTONDOWN:
20         g_clicks_xy.append((x, y))
21
22 """
23 @im: The input image as read (in BGR)
24 """
25 def get_model(im):
26     """
27     Process user input (clicks) and Extract the target image
28     (model) in hsv.
29     """
30     global g_clicks_xy
31     hsv = cv2.cvtColor(im, cv2.COLOR_BGR2HSV)
32     cv2.namedWindow('2 clicks, then press q')
33     cv2.setMouseCallback('2 clicks, then press q', on_click)
34     while True:
35         cv2.imshow('2 clicks, then press q', im)
36         k = cv2.waitKey()
37         if k == ord('q'):
38             break
39     cv2.destroyAllWindows()
40     # for clicks, index 0 = x, index 1 = y
41     g_clicks_xy = sorted(g_clicks_xy,
42                          key = lambda x: x[0]**2 + x[1]**2,
43                          reverse = True) [-2:]
44     click_br, click_tl = g_clicks_xy[0], g_clicks_xy[1]
45     w = click_br[0] - click_tl[0]
46     h = click_br[1] - click_tl[1]
47     hsvt = hsv[click_tl[1]: click_tl[1] + h,
48              click_tl[0]: click_tl[0] + w]
49     rect = cv2.rectangle(im.copy(), g_clicks_xy[0], g_clicks_xy[1],
50                          color = (0, 255, 0))
51     qimshow(rect, wname = 'selection', delay = 3)
52     return hsvt
53
54 def main():
55     cap = cv2.VideoCapture(sys.argv[1])
56     # take first frame of the video
57     ret, frame = cap.read()
58     # setup initial location of window
59     r, h, c, w = 250, 90, 400, 125 # simply hardcoded the values
60     track_window = (c, r, w, h)
61     # set up the ROI for tracking
```

```

62     roi = frame[r:r+h, c:c+w]
63     hsv_roi = get_model(frame)
64     click_br, click_tl = g_clicks_xy[0], g_clicks_xy[1]
65     w = click_br[0] - click_tl[0]
66     h = click_br[1] - click_tl[1]
67     x = click_tl[0]
68     y = click_tl[1]
69     track_window = (x, y, w, h)
70     #hsv_roi = cv2.cvtColor(roi, cv2.COLOR_BGR2HSV)
71     # Filter out noise
72     mask = cv2.inRange(hsv_roi, np.array((0., 60.,32.)),
73                        np.array((180.,255.,255.)))
74     roi_hist = cv2.calcHist([hsv_roi],[0],mask,[180],[0,180])
75     cv2.normalize(roi_hist,roi_hist,0,255,cv2.NORM_MINMAX)
76     # Setup the termination criteria,
77     #either 10 iteration or move by atleast 1 pt
78     term_crit = ( cv2.TERM_CRITERIA_EPS | cv2.TERM_CRITERIA_COUNT, 10, 1 )
79     while(1):
80         ret ,frame = cap.read()
81         if ret == True:
82             hsv = cv2.cvtColor(frame, cv2.COLOR_BGR2HSV)
83             dst = cv2.calcBackProject([hsv],[0],roi_hist,[0,180],1)
84             # apply meanshift to get the new location
85             ret, track_window = cv2.meanShift(dst, track_window, term_crit)
86             # Draw it on image
87             x,y,w,h = track_window
88             img2 = cv2.rectangle(frame, (x,y), (x+w,y+h), 255,2)
89             cv2.imshow('img2',img2)
90             k = cv2.waitKey(30) & 0xff
91             if k == ord('q'):
92                 break
93         else:
94             break
95     cv2.destroyAllWindows()
96     cap.release()
97
98 main()

```