

Computer Vision

James L. Crowley

M2R MoSIG
GVR and UIS

Fall Semester
15 Oct 2015

Lesson 3

Bayesian Detection and Tracking of Objects

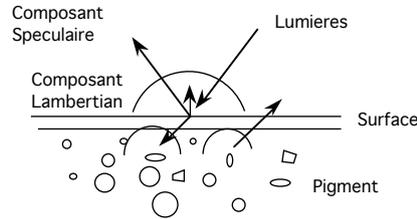
Lesson Outline:

1	Detection and Tracking using Color	2
1.1	Separating Specular and Lambertian Reflection.....	2
1.2	Histograms as an Estimation of Probability of Color.....	3
1.3	Bayesian Object Detection with Color	4
1.4	Color Skin Detection.....	6
2	Bayesian Tracking of Gaussian Blobs	8
2.1	Moment Calculations for Blobs.....	8
2.2	Bayesian Tracking	12
2.3	Temporal Prediction.....	13
2.4	Detecting the target.....	16
2.5	Updating the Estimated Blob Parameters	16
2.6	Managing Lost Targets	17
3	The Kalman Filter	18
3.1	State Vector	19
3.2	Confidence and Uncertainty.....	19
3.3	State Estimation	22
3.4	Summary of Kalman filter equations.....	22

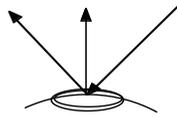
1 Detection and Tracking using Color

Recall the Bi-Chromatic reflection function:

$$R(i, e, g, \lambda) = \alpha R_s(i, e, g, \lambda) + (1 - \alpha) R_L(i, \lambda)$$



For Lambertian reflection, the intensity is generally determined by surface orientation, while color is determined by Pigment.



1.1 Separating Specular and Lambertian Reflection.

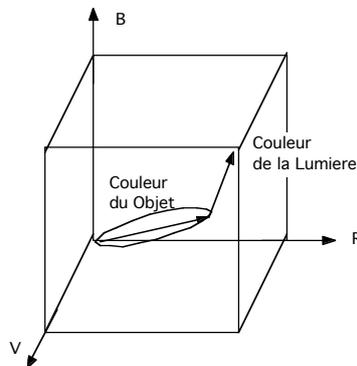
Consider what happens at a specular reflection.

The specularity has the same spectrum as the illumination.

The rest of the object has a spectrum that is the product of illumination and pigments. Suppose that the image be composed of color pixels, $\vec{c}(i, j)$

We can construct an RGB color histogram using a 3 Dimensional table $h(R, G, B)$.

$$\forall \vec{c}(i, j) : h(\vec{c}(i, j)) = h(\vec{c}(i, j)) + 1$$



Suppose the image contains only a single (non-planar) object that obeys a bi-chromatic reflection model. If the image contains a specular reflection, then two clear axes will emerge.

One axis from the origin to the RGB of the product of the illumination and the source. The other axis towards the RGB representing the illumination.

We can fit lines to these two axes. The dominant line indicates the object pigment color modified by the illumination; $P(\lambda)S(\lambda)$. The line intersects the edge of the cube at this color.

The second (less dominant line) indicates the source color $S(\lambda)$. Source color can be observed from the point where this line exits the cube.

Projecting ALL pixels onto these two lines gives two images:

- 1) A Lambertian image (without the specularly) and
- 2) A specular image (without the Lambertian component).

Lesson: color statistics, particularly histograms, can provide powerful analysis tools .

1.2 Histograms as an Estimation of Probability of Color

A histogram is a table of frequency of occurrence. Histograms have many uses in image processing and computer vision. One such use is for pixel level probabilistic detection of objects based on local appearance. A simple example is the use of color histograms to detect objects based on color statistics.

Assume that we have a color image, where each pixel (i,j) is a color vector, $\vec{c}(i,j)$, composed of 3 integers between 0 and 255 representing Red, Green and Blue.

$$\vec{c}(i,j) = \begin{pmatrix} R \\ G \\ B \end{pmatrix}$$

We can build a color histogram of the image by counting the number of times that each unique value of (R, G, B) occurs in the image. To do this we allocate a table $h(r, g, b)$ of $256 \times 256 \times 256$ cells, with each cell initially set to zero. The table $h(r,g,b)$ has $256^3 = 2^{24} = 16$ Mega Pixels.

We then visit each pixel and add one to the value of the cell that corresponds to its value of R, G, B

$$\forall_{i,j} h(\vec{c}(i,j)) = h(\vec{c}(i,j)) + 1$$

The table $h(\vec{c})$ then represents the frequency of occurrence for each possible color vector \vec{c} . Given that the image is composed of M pixels, then this table tells us the probability of finding a pixel of color \vec{c} at any position in the image.

$$P(\vec{c}) = \frac{1}{M} h(\vec{c})$$

If we take a second image of the same scene, and we assume that the color and illumination are similar, we can use the histogram to predict the probability of color vectors in the second image.

1.3 Bayesian Object Detection with Color

Suppose that we have K RGB images composed of $R \cdot C$ pixels, $C_k(i, j)$. This gives a total of $M = K \cdot R \cdot C$ pixels. Suppose that we have a subset, T (for target) of these pixels that belong to a target class. Suppose that T contains M_T pixels.

We allocate two tables $h(r, g, b)$ and $h_t(r, g, b)$ and as before and use these to construct two histograms.

$$\forall_{i,j,k} h(\vec{c}_k(i,j)) = h(\vec{c}_k(i,j)) + 1$$

$$\forall_{(i,j,k) \in T} h_T(\vec{c}_k(i,j)) = h_T(\vec{c}_k(i,j)) + 1$$

For any color vector, \vec{c} , have TWO probabilities :

$$P(\vec{c}) = \frac{1}{M} h(\vec{c}) \quad \text{and} \quad P(\vec{c} | T) = \frac{1}{M_T} h_T(\vec{c})$$

Bayes rule tells us that we can estimate the probability that a pixel belongs to an object given its color as:

$$P(T | \vec{c}) = \frac{P(\vec{c} | T)P(T)}{P(\vec{c})}$$

We have $P(\vec{c} | T)$ and $P(\vec{c})$. $P(T)$ is the probability that a pixel belongs to a target. For the training image, this is given by

$$P(T) = \frac{M_T}{M}$$

From this we can show that the probability of a target, T, is simply the ratio of the two tables.

$$P(T | \bar{c}) = \frac{P(\bar{c} | T)P(T)}{P(\bar{c})} = \frac{\frac{1}{M_T} h_t(\bar{c}) \cdot \frac{M_t}{M}}{\frac{1}{M} h(\bar{c})} = \frac{h_T(\bar{c})}{h(\bar{c})}$$

We can use this to compute a lookup table $L_T(\bar{c})$: $L_T(\bar{c}) = \frac{h_T(\bar{c})}{h(\bar{c})}$

If we ASSUME that a new image, $x(i,j)$, has similar illumination and color composition then we can use this technique to assign a probability to each pixel by table lookup. The result is an image in which each pixel is a probability $T(i,j)$ that the pixel (i,j) belongs to the target T.

$$T(i,j) = L_T(x(i,j))$$

The reliability is improved by using more training images.

The naive statistics view says to have at least 8 training samples for histogram cell. For example, in our RGB example, $h(c)$ was composed of

$$Q = 2^8 \cdot 2^8 \cdot 2^8 = 2^{24} \text{ cells.}$$

Thus we need $M = 2^3 \cdot 2^{24} = 2^{27}$ training pixels. This is not a problem for $P(\bar{c}) = \frac{1}{M} h(\bar{c})$ but may be a problem for $P(\bar{c} | T) = \frac{1}{M_T} h_T(\bar{c})$.

(Note that a 1024 x 1024 image contains 2^{20} pixels. This is the definition of 1 Mega)

Q is the "capacity" of the histogram, measured as the number of cells.

$Q = N^D$ where N is the number of values per feature and D is the number of features.

A more realistic view is that the training data must contain a variety of training samples that reflect that variations in the real world.

What can we do? Often we can reduce both the number, D, of features and the number of values, N, for each feature.

For example, for many color images, N=32 color values are sufficient to detect objects. We simply divide each color R, G, B by 8.

$$R' = \text{Trunc}(R/8), \quad G' = \text{Trunc}(G/8), \quad B' = \text{Trunc}(B/8).$$

We can also use our knowledge of physics to look for features that are "invariant".

1.4 Color Skin Detection

Luminance captures surface orientation (3D shape) while Chrominance is a signature for object pigment (identity). Thus it is convenient to transform the (RGB) color pixels into a color space that separates Luminance from Chrominance.

$$\begin{pmatrix} L \\ C_1 \\ C_2 \end{pmatrix} \Leftarrow \begin{pmatrix} R \\ G \\ B \end{pmatrix}$$

Normalizing out luminance provides a popular space for skin detection: the (r,g) space.

Luminance: $L = R + G + B$

$$\text{Chrominance : } r = c_1 = \frac{R}{R + G + B} \quad g = c_2 = \frac{G}{R + G + B}$$

These are often called "r" and "g" in the literature. The (r, g) space is often used to detect skin colored pixels. It is common to normalize r and g to natural numbers coded with N values between 0 and N – 1 by :

$$r = \text{trunc}\left(N \cdot \frac{R}{R + G + B}\right) \quad g = \text{trunc}\left(N \cdot \frac{G}{R + G + B}\right)$$

Skin pigment is generally always the same chrominance value. Luminance can change with pigment density, and skin surface orientation, but chrominance will remain invariant.

Thus we can use $\begin{pmatrix} r \\ g \end{pmatrix}$ as an invariant color signature for detecting skin in images.

Suppose we have a set of K training images $\{c_k(i,j)\}$ of size RxC where each pixel is an RGB color vector. This gives a total of $M = K \times R \times C$ color pixels. Suppose that M_{skin} of these are labeled as skin pixels.

We can simplify our technique by projecting these onto chrominance pixels. From experience, $N = 32$ color values seems to work well for skin.

We allocate two table : $h(r,g)$ and $h_{skin}(r,g)$ of size $N \times N$.

For all i,j,k in the training set $\{c_k(i,j)\}$:

BEGIN

$$r = \text{trunc}\left(N \cdot \frac{R}{R+G+B}\right) \quad g = \text{trunc}\left(N \cdot \frac{G}{R+G+B}\right)$$

$$h(r,g) = h(r,g) + 1$$

IF the pixel $c_k(i,j)$ is skin THEN

$$h_{\text{skin}}(r,g) = h_{\text{skin}}(r,g) + 1$$

END

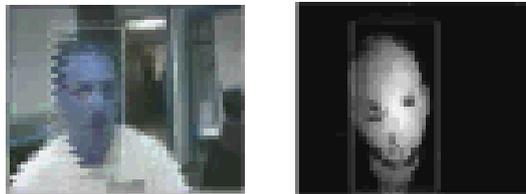
As before, we can obtain a lookup table $L_{\text{skin}}(r,g)$ that gives the probability that a pixel is skin.

$$L_{\text{skin}}(r,g) = \frac{h_{\text{skin}}(r,g)}{h(r,g)}$$

Given a new RGB image $C(i,j)$:

$$r = \text{trunc}\left(N \cdot \frac{R}{R+G+B}\right) \quad g = \text{trunc}\left(N \cdot \frac{G}{R+G+B}\right)$$

$$T_{\text{skin}}(i,j) = L_{\text{skin}}(r(i,j), g(i,j))$$



(images from a Bayesian skin tracking in real time - 1993)

We can improve the detection by tracking skin colored regions.

2 Bayesian Tracking of Gaussian Blobs

Rather than represent a skin region as a collection of pixels, we can calculate a Gaussian Blob. A "Blob" represents a region of an image. Gaussian blobs express a region in terms of moments.

Assume of image of probabilities of the detection of a target: $T(i,j)$, where for each pixel:

$$T(i,j) = L_T(r(i,j), g(i,j))$$

The zeroth moment of the probabilities is the mass (sum of probabilities). Average mass represents confidence.

The first moment gives is the center of gravity. This is the "position" of the blob.

The second moment is the covariance. This gives size and orientation.

We typically enclose the blob in some rectangular Region of Interest (ROI) in order to avoid "distraction" by neighboring blobs. The ROI is obtained by some form of estimation or a priori knowledge. In continuous operation the ROI be provided by tracking.

Let us represent the ROI as a rectangle : (t,l,b,r)

- t - "top" - first row of the ROI.
- l - "left" - first column of the ROI.
- b - "bottom" - last row of the ROI
- r - "right" -last column of the ROI.

(t,l,b,r) can be seen as a bounding box, expressed by opposite corners (l,t), (r,b)
We will compute the moments within this ROI (bounding box).

2.1 Moment Calculations for Blobs

Given a target probability image $T(i,j)$ and a ROI (t,l,b,r):

$$\text{Sum: } S = \sum_{i=l}^r \sum_{j=t}^b T(i,j)$$

We can estimate the "confidence" as the average detection probability:

Confidence: $CF = \frac{S}{(b-t)(r-l)}$

First moments:

$$x = \mu_i = \frac{1}{S} \sum_{i=l}^r \sum_{j=t}^b T(i,j) \cdot i$$

$$y = \mu_j = \frac{1}{S} \sum_{i=l}^r \sum_{j=t}^b T(i,j) \cdot j$$

Position is the center of gravity: (μ_i, μ_j)

We will use this as the position of the blob.

Second Moments:

$$\sigma_i^2 = \frac{1}{S} \sum_{i=l}^r \sum_{j=t}^b T(i,j) \cdot (i - \mu_i)^2$$

$$\sigma_j^2 = \frac{1}{S} \sum_{i=l}^r \sum_{j=t}^b T(i,j) \cdot (j - \mu_j)^2$$

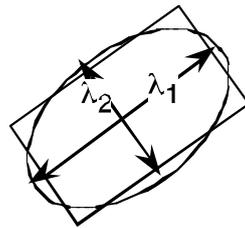
$$\sigma_{ij}^2 = \frac{1}{S} \sum_{i=l}^r \sum_{j=t}^b T(i,j) \cdot (i - \mu_i) \cdot (j - \mu_j)$$

These compose a covariance matrix: $C = \begin{pmatrix} \sigma_i^2 & \sigma_{ij}^2 \\ \sigma_{ij}^2 & \sigma_j^2 \end{pmatrix}$

The principle components (λ_1, λ_2) determine the length and width.

The principle direction determines the orientation of the length.

We can discover these by principle components analysis.



$$RCR^T = \Lambda = \begin{pmatrix} \lambda_1^2 & 0 \\ 0 & \lambda_2^2 \end{pmatrix}$$

where

$$R = \begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix}$$

The length to width ratio, λ_1/λ_2 , is an invariant for shape.

The angle θ is a ‘‘Covariant’’ for orientation.

We can use the “eigenvalues”, or characteristic values, λ_1, λ_2 , to define the “width and height” of the blob:

For example:

$$w=\lambda_1, h=\lambda_2$$

This suggests a "feature vector" for the blob: $\vec{X} = \begin{pmatrix} x \\ y \\ w \\ h \\ \theta \end{pmatrix}$

where $x = \mu_i, y = \mu_j, w = \lambda_1, h = \lambda_2$ and $CF = \frac{S}{(b-t)(r-l)}$

The confidence (CF) can be seen as the “Likelihood” that the model for the blob is correct.

Tracking allows us to continually update an estimate for the features of the Blob, even if the blob is temporarily lost to occlusion or noise.

The tracked object is often referred to as a "target". The vector \vec{X} provides the “model” for the target blob:

$$\text{Blob model: } \vec{X} = \begin{pmatrix} x \\ y \\ w \\ h \\ \theta \end{pmatrix} \quad \text{Precision: } P = \begin{pmatrix} \sigma_{xx}^2 & \sigma_{xy}^2 & \sigma_{xw}^2 & \sigma_{xh}^2 & \sigma_{x\theta}^2 \\ \sigma_{yx}^2 & \sigma_{yy}^2 & \sigma_{yw}^2 & \sigma_{yh}^2 & \sigma_{y\theta}^2 \\ \sigma_{wx}^2 & \sigma_{wy}^2 & \sigma_{ww}^2 & \sigma_{wh}^2 & \sigma_{w\theta}^2 \\ \sigma_{hx}^2 & \sigma_{hy}^2 & \sigma_{hw}^2 & \sigma_{hh}^2 & \sigma_{h\theta}^2 \\ \sigma_{\theta x}^2 & \sigma_{\theta y}^2 & \sigma_{\theta w}^2 & \sigma_{\theta h}^2 & \sigma_{\theta\theta}^2 \end{pmatrix}$$

along with CF. (Confidence)

The covariance P expresses our “uncertainty” about the values of each of the parameters. This is also called the “Precision”.

Formally, precision is defined as the 2nd moment of the error of the blob model.

Let \vec{X} be the TRUE (unknown) model and $\hat{\vec{X}}$ be the estimated value. The error, \vec{E} , is the difference

$$\vec{E} = \vec{X} - \hat{\vec{X}}$$

For a set of M target observations, the average error is the expected value of \vec{E}

$$\bar{\mu}_E = E\{\bar{E}\} = \frac{1}{M} \sum_{m=1}^M \bar{E}$$

This is also known as a “bias” in the model.
The precision is the covariance of the error.

$$P = E\{(\bar{E} - \bar{\mu}_E)(\bar{E} - \bar{\mu}_E)^T\} = \frac{1}{M} \sum_{m=1}^M (\bar{E} - \bar{\mu}_E)(\bar{E} - \bar{\mu}_E)^T$$

This is a covariance matrix for the model components.

$$P = \begin{pmatrix} \sigma_{xx}^2 & \sigma_{xy}^2 & \sigma_{xw}^2 & \sigma_{xh}^2 & \sigma_{x\theta}^2 \\ \sigma_{yx}^2 & \sigma_{yy}^2 & \sigma_{yw}^2 & \sigma_{yh}^2 & \sigma_{y\theta}^2 \\ \sigma_{wx}^2 & \sigma_{wy}^2 & \sigma_{ww}^2 & \sigma_{wh}^2 & \sigma_{w\theta}^2 \\ \sigma_{hx}^2 & \sigma_{hy}^2 & \sigma_{hw}^2 & \sigma_{hh}^2 & \sigma_{h\theta}^2 \\ \sigma_{\theta x}^2 & \sigma_{\theta y}^2 & \sigma_{\theta w}^2 & \sigma_{\theta h}^2 & \sigma_{\theta\theta}^2 \end{pmatrix}$$

In most cases, the true value is unknown, and so the precision must be estimated from the difference between the predicted and observed blobs.

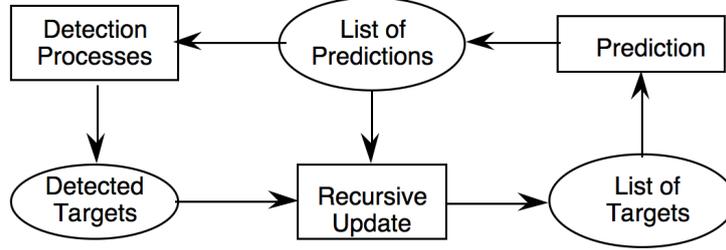
This precision can be initialized to a very small value. Our uncertainty about the precision will grow over time between observations, due to unobserved motion, acceleration, etc. We can estimate this growth by the difference between a predicted blob and its observation. This process will tend to reveal covariance between terms of the matrix.

If possible it is preferred to estimate P from some known values of \vec{X} or from knowledge about the sensors and from “experience”.

In the absence of information about the correlation between features in X , the Covariance can be initialized as a diagonal matrix. The tracking process will review any correlations.

2.2 Bayesian Tracking

A Bayesian tracker is a recursive estimator, composed of three phases: Predict, Detect, Update.



Detection can be provided by detecting the blob using color statistics within a target "Region of Interest" given by a bounding box centered on a previous position. The size of this box is determined by the estimated size of the blob enlarged by the uncertainty P .

The Gaussian window is the previous covariance for the blob, enlarged by some "uncertainty" covariance. The uncertainty captures the possible loss of information during the time from the most recent observation.

Our Gaussian blob is

Position: $\bar{\mu}_t = \begin{pmatrix} \mu_i \\ \mu_j \end{pmatrix}$ Size : $C_t = \begin{pmatrix} \sigma_i^2 & \sigma_{ij}^2 \\ \sigma_{ij}^2 & \sigma_j^2 \end{pmatrix}$ along with CF_t .

where the second moment of the detected pixels, C , was used to compute to determine the width, height and orientation:

$$RCR^T = \Lambda = \begin{pmatrix} \lambda_1^2 & 0 \\ 0 & \lambda_2^2 \end{pmatrix} = \begin{pmatrix} h^2 & 0 \\ 0 & w^2 \end{pmatrix}$$

so

$$C = R^T \Lambda R = \begin{pmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{pmatrix} \begin{pmatrix} h^2 & 0 \\ 0 & w^2 \end{pmatrix} \begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix}$$

$$R = \begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix}$$

Let us define the estimated blob at time t as: $\hat{\mu}_t, \hat{C}_t$

Let us define the predicted feature vector at time t as: $\bar{\mu}_t^*, C_t^*$

The following describes prediction and estimation (updating).

In the absence of any knowledge of movement, we can predict that the blob will be at the last observed position. This is written as :

$$\bar{X}_t^* \leftarrow \hat{X}_{t-1}$$

This is called a "process model". The process model predicts the state vector at time t given the estimate at time t-Δt:

$$X_t^* = \operatorname{argmax} \{ P(X_t^* | X_{t-\Delta t}) \}$$

This can be written as:

$$\begin{aligned} X_t^* &:= \vartheta(\Delta t) \hat{X}_{t-\Delta t} + R \\ P_t^* &:= \vartheta(\Delta t) \hat{P}_{t-\Delta t} \vartheta(\Delta t)^T + Q_x \end{aligned}$$

For the simple zeroth order filter $\vartheta(\Delta t)$ is the identity matrix.

$$\vartheta(\Delta t) = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

For higher order filters, the state vector contains temporal derivatives and $\vartheta(\Delta t)$ permits using the derivatives to predict future state $\bar{X}_t^* \leftarrow \bar{X}_{t-1} + \Delta t \bar{X}'_{t-1}$ where X' is the first temporal derivative.

R is the expected "error" of the process model, due to un-modeled derivatives and noise. The expectation of the error of this prediction and is generally 0.

$$R = E\{\bar{X} - \hat{X}\} = 0$$

because the prediction is zero mean error.

Q is the second moment of the error. It is NOT zero.

$$Q = E\{(\bar{X} - \hat{X})(\bar{X} - \hat{X})^T\}$$

Q captures the loss of precision in the evolution of the process noise.

In a first order Kalman filter this would have been

$$\vec{X}_t^* \leftarrow \vec{X}_{t-1} + \Delta t \vec{X}'_{t-1}$$

where

$$\vec{X}'_{t-1} = \frac{d}{dt} \begin{pmatrix} x \\ y \\ w \\ h \\ \theta \end{pmatrix} = \begin{pmatrix} x' \\ y' \\ w' \\ h' \\ \theta' \end{pmatrix}$$

and Δt is the time step. But let us keep the explanation simple for now.

Because of unobserved velocity and acceleration, etc, the estimate of the blob grows more uncertain with time. This is estimated by a prediction error : $Q_{\Delta t}$.

$$P_t^* = \hat{P}_{t-\Delta t} + Q_{\Delta t} \Delta t^2$$

The values in $Q_{\Delta t}$ can be "calibrated" by measuring the average error between predicted and observed blobs from a labeled training sequence, or it can be "estimated".

We then use the predicted target blob to compute a new predicted ROI for the detection as described above.

2.4 Detecting the target

For each sensor, a predicted sensor signal Y_t^* is generated based on current the estimated system state X_t^* .

$$Y_t^* = \operatorname{argmax}\{p(Y_t | X_t^*)\}.$$

We obtain this prediction from the predicted target model.

$$Y_t^* = H_X^Y X_t^*$$

In the trivial case of a zeroth order tracker, if \bar{X} and \bar{Y} have the same features than H_X^Y is the identity matrix.

$$H_X^Y = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

Else H_X^Y can "extract" observed features from a more complex model.

In the case of blob detection, we can estimate the bounding box using all 5 parameters of Y_t^* . In this case H_X^Y is the identity matrix.

2.5 Updating the Estimated Blob Parameters

Detection can introduce new errors, such distraction by adjacent target

To minimize such errors, we weight the detected target by a predicted uncertainty as was shown above. This is a form of robust estimation that rejects outlying detections.

We weight the detected pixels by the predicted Normal density.

$$T(i, j) \leftarrow L(\bar{c}(i, j)) \cdot e^{-\frac{1}{2} \begin{pmatrix} i \\ j \end{pmatrix} - \begin{pmatrix} \mu_i \\ \mu_j \end{pmatrix}}^T 2C_t^{-1} \begin{pmatrix} i \\ j \end{pmatrix} - \begin{pmatrix} \mu_i \\ \mu_j \end{pmatrix}}$$

Where $T(i, j)$ is the new observed probability image at the new time, t.

and $L()$ is our lookup table for the ration of histograms as explained above.

We then estimate the new position and covariance using this product:

$$\text{First moments:} \quad \hat{\mu}_i = \frac{1}{S} \sum_{i=l}^r \sum_{j=t}^b T(i, j) \cdot i \quad \hat{\mu}_j = \frac{1}{S} \sum_{i=l}^r \sum_{j=t}^b T(i, j) \cdot j$$

Second Moments:

$$\hat{\sigma}_i^2 = \frac{1}{S} \sum_{i=1}^r \sum_{j=1}^b T(i, j) \cdot (i - \hat{\mu}_i)^2$$

$$\hat{\sigma}_j^2 = \frac{1}{S} \sum_{i=1}^r \sum_{j=1}^b T(i, j) \cdot (j - \hat{\mu}_j)^2$$

$$\hat{\sigma}_{ij}^2 = \frac{1}{S} \sum_{i=1}^r \sum_{j=1}^b T(i, j) \cdot (i - \hat{\mu}_i) \cdot (j - \hat{\mu}_j)$$

which gives: $\hat{\mu}_t = \begin{pmatrix} \hat{\mu}_i \\ \hat{\mu}_j \end{pmatrix}$ and $\hat{C}_t = \begin{pmatrix} \hat{\sigma}_i^2 & \hat{\sigma}_{ij}^2 \\ \hat{\sigma}_{ij}^2 & \hat{\sigma}_j^2 \end{pmatrix}$

From which we can update the new estimate : \hat{X}_t

2.6 Managing Lost Targets

Targets can disappear due to occlusion or lost tracking. For stability we accumulate confidence of targets over time.

$$CF_t = \alpha CF + (1 - \alpha) CF_{t-\Delta t} +$$

CF_{\min} is the minimum required average probability per pixel to detect a target.

If $CF_t \leq CF_{\min}$ then a target is removed.

The weight α determines the decay rate for lost targets.

3 The Kalman Filter

The Kalman filter is a form of Bayesian estimator. Unlike the Bayesian estimation process above, the Kalman filter updates the estimated state vector from the difference of predicted and observed sensor data. (rather than using robust estimation).

Three key steps characterise Bayesian estimation problems (including Kalman filters).

1) A process model: The process model predicts the state vector at time t given the estimate at time $t-\Delta t$:

$$X_t^* = \operatorname{argmax} \{ p(X_t^* | X_{t-\Delta t}) \}$$

2) A sensor model: For each sensor, a predicted sensor signal Y_t^* is generated based on current the estimated system state X_t^* .

$$Y_t^* = \operatorname{argmax} \{ p(Y_t | X_t^*) \}.$$

3) Re-estimation: A new estimated value, X_t is computed based on information provided by the difference between the predicted and observed sensor values.

$$X_t = \operatorname{argmax} \{ p(X_t | X_t^*, Y_t - Y_t^*) \}$$

The Kalman filter uses a linear dynamic model to provide these estimates. That is, the process model and sensor models are represented by linear equations. A fixed time step and previously estimated derivative values are used to estimate the current value of the state variables. A quadratic form this same dynamic equation is used to predict the error of the state vector.

A simple zeroth order Kalman filter may be used to track bodies, faces and hands in video sequences. In the following example, let us assume a target whose properties are represented by a "state vector" composed of position, width, height and orientation (x, y, w, h, θ) .

A 5x5 covariance matrix is associated with this vector to represent correlations in errors between parameters. Although prediction does not change the estimated position, it does enlarge the uncertainties of the position and size of the expected target.

3.1 State Vector

The target state vector, \hat{X}_t , is composed of the position, scale and orientation of the target. For example, this can represent a human face.

$$\bar{X} = \begin{pmatrix} x \\ y \\ w \\ h \\ \theta \end{pmatrix}$$

where x, y are the position of the target in pixels
 w, h are the width and height of the target, and
 θ is the image plane orientation of the target.

In the case of a first order filter, each of the parameters is accompanied by first temporal derivative.

$$\hat{X}_t = \begin{pmatrix} x \\ \dot{x} \\ y \\ \dot{y} \\ w \\ \dot{w} \\ h \\ \dot{h} \\ \theta \\ \dot{\theta} \end{pmatrix}$$

Where the dot indicates temporal derivative.

$$\dot{x} = \frac{dx}{dt}$$

The Kalman filter equations are able to use information from the difference of observed and predicted state to estimate the temporal derivatives.

3.2 Confidence and Uncertainty

The state vector is accompanied by a covariance matrix and a confidence factor. The confidence factor is an integer between 0, and a maximum confidence value.

$$CF_t \in [0, CF_{\max}]$$

The position uncertainty (or precision) is the covariance matrix for the state vector

For the case of a first order filter, this matrix becomes 10 by 10 with covariance between all terms.

For each target, at each time, t , the tracker maintains an estimated state \hat{X}_t , as well as its precision \hat{P}_t and confidence factor, CF_t . Based on a previous state and precision \hat{X}_{t-1} , \hat{P}_{t-1} and CF_{t-1} as well as the observation \hat{P}_t from detection function accompanied by the observed precision P_y , and detection confidence CF_y .

$$\hat{X}_t, \hat{P}_t, CF_t = F\{X_{t+\Delta t}^*, P_{t+\Delta t}^*, CF_{t-\Delta t}, Y, P_y, CF_y\}$$

Given a target at time $t-\Delta t$, the prediction equations predict its new position, and validation gate at time t . The general form of the prediction equations are :

$$\begin{aligned} X_t^* &:= \vartheta(\Delta t)\hat{X}_{t-\Delta t} + R \\ P_t^* &:= \vartheta(\Delta t)\hat{P}_{t-\Delta t}\vartheta(\Delta t)^T + Q_x \end{aligned}$$

where R is the expected error or “residue” of the temporal prediction (generally 0) and Q is its second moment.

These equations are a linear estimation of movement based on a Taylor series approximation.

The term R is a residue that represents higher order (non-estimated) derivatives. This expected value of this term is zero and thus R is commonly omitted. The second moment of R represents the uncertainty due to accelerations (and higher order derivatives). Thus second moments, Q , estimates the loss of precision due to higher order terms.

$$Q = E \{R R^T\}$$

When included in the prediction, the term Q provides to an additive growth in the validation gate that is translated to the search region. When a target is detected, this growth disappears in the update phase. However if no target is detected, the result is that the validation gate (and thus the region of interest) grows with each frame until the target is re-acquired or until the target is declared lost.

For a first order filter, the prediction matrix $\vartheta(\Delta t)$ predicts new values as a function of the time step and the estimated derivatives.

$$\Phi(\Delta t) = \begin{pmatrix} 1 & \Delta t & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & \Delta t \\ 0 & 0 & \cdots & 0 & 1 \end{pmatrix}$$

In the case of face tracking, we generally assume that face accelerations are too rapid to be estimated. Thus we may estimate only target position (order 0 tracking) or position and velocity (order 1 tracking). In this case, the prediction matrix, $\Phi(\Delta t)$, becomes a trivial identity matrix and our prediction equations are reduced to

$$\begin{aligned} X_t^* &= \hat{X}_{t-\Delta t} \\ P_t^* &= \hat{P}_{t-\Delta t} + Q_x \end{aligned}$$

A validation gate provides a region of interest (ROI), which serves to limit estimation to the region where the target can be detected. For face detection, this region of interest specifies a range of positions and scales and possibly orientations at which the target may be sought. This greatly accelerates processing by avoiding processing unnecessary pixels.

$$\text{ROI} = (x_{\min}, x_{\max}, y_{\min}, y_{\max}, w_{\min}, w_{\max}, h_{\min}, h_{\max}, \theta_{\min}, \theta_{\max})$$

The ROI is computed using the scale to define the width and height of a rectangular region over which the target will be sought. A ROI based on three standard deviations reasonable size search region. Such a ROI is defined can be defined by adding 2 or 3 standard deviations to the precision.

This can then be used to drive face detection process using a cascade of linear classifiers. The resulting face location is noted as the observed location

$$Y_t = \begin{pmatrix} x \\ y \\ w \\ h \\ \theta \end{pmatrix}$$

The difference between the predicted and observed face locations is known as the innovation:

$$\text{Innovation: } (Y_t - H_x^y X_t^*) = (Y_t - Y_t^*)$$

(where H is the identity matrix for our example of face detection.)

3.3 State Estimation

This scale is used to update the precision of the face position estimate. The estimated state vector is updated by using the difference between predicted and estimated position as an "innovation". In the case of the 0th order filter, the equations are relatively simple. We first compute a Kalman "Gain matrix:

$$K = P_t^* (P_t^* - P_y)^{-1}$$

from this we can compute the following update equations:

$$\begin{aligned}\hat{X}_t &= X_t^* + K(Y_t - H_X^Y X_t^*) \\ \hat{P}_t &= P_t^* - K P_t^*\end{aligned}$$

The resulting algorithm in the case of a 1st order filter, a transformation matrix, ϑ : used to project X onto the observed variables Y.

3.4 Summary of Kalman filter equations.

1) Process model: Temporal prediction (using linear dynamic models).

$$\begin{aligned}X_t^* &:= \vartheta(\Delta t) \hat{X}_{t-\Delta t} + R && \text{Temporal prediction of State} \\ P_t^* &:= \vartheta(\Delta t) \hat{P}_{t-\Delta t} \vartheta(\Delta t)^T + Q_x && \text{Precision of prediction}\end{aligned}$$

2) Sensor Model: Observation prediction - H_X^Y predicts observed state from estimated state

$$\begin{aligned}Y_t^* &= H_X^Y X_t^* && \text{Predicted observation} \\ P_y &= H_X^Y P_t^* H_X^Y && \text{Precision of prediction}\end{aligned}$$

and

$$\begin{aligned}Y_t & \text{ actual observation at time } t \\ P_y &= E\{WW^T\} \text{ Calibrated precision of estimation}\end{aligned}$$

3) Observation and update:

$$\begin{aligned}K &= P_t^* H_X^{YT} (H_X^Y P_t^* H_X^{YT} - P_y)^{-1} && \text{(Kalman Gain)} \\ \hat{X}_t &= X_t^* + K(Y_t - H_X^Y X_t^*) && \text{Update of estimated state} \\ \hat{P}_t &= P_t^* - K H_X^Y P_t^* && \text{Update precision of estimated state}\end{aligned}$$