

Pattern Recognition and Machine Learning

James L. Crowley

ENSIMAG 3 - MMIS
Lesson 3

Fall Semester 2017
25 October 2017

Detecting and Locating Faces with Color: Three Challenges for Project 1.

Outline

Notation	2
1. Lab 1: Face Detection using Color.....	3
1.1 Algorithm Overview	3
2. Challenge 1: Detecting skin pixels with color	4
2.1 Baseline: RGB histograms	4
2.2 Variation: Detection with Chrominance (normalized color). 7	
2.3 Free Parameters	8
3. Challenge 2: Detecting Faces with Color.....	9
3.2 Baseline: Sliding Window Detector.....	10
3.3 Variation: Robust Detection using a Gaussian mask.	12
3.4 Free parameters to test	13
4. Challenge 3: Face Localization	14
4.1 Detection vs Localization (Parameter Estimation)	14
4.2 Baseline: Localization by non-maximum suppression	15
4.3 Variation: Localization by Robust Estimation	15
4.4 Free parameters to test	16

Notation

$C(i,j)$	An RGB image of size $I \times J$ pixels, 8 bits per color
$\{C_k(i,j)\}$	A set of K images for training.
\bar{X}	An observation.
$\{\bar{X}_m\}$	Training data for learning.
$y(\bar{X}_m)$	An annotation (or ground truth) function $y(\bar{X}_m) \in \{P, N\}$
$g(\bar{X}_m)$	Discriminant function. $0 \leq g(\bar{X}_m) \leq 1$
M	The number of training samples.
M_T	The number of training samples in the target class
$h(\bar{X})$	A multidimensional histogram of for \bar{X}
Q	The number of cells in a histogram

$$G(i, j; \bar{\mu}, \Sigma) = \frac{1}{2\pi \det(\Sigma)^{1/2}} e^{-\frac{1}{2}(\bar{p}-\bar{\mu})^T \Sigma^{-1}(\bar{p}-\bar{\mu})}$$

A Gaussian mask.

$$\bar{\mu} = \begin{pmatrix} \mu_i \\ \mu_j \end{pmatrix}$$

the center position

$$\Sigma = \begin{pmatrix} \sigma_i^2 & \sigma_{ij} \\ \sigma_{ij} & \sigma_j^2 \end{pmatrix}$$

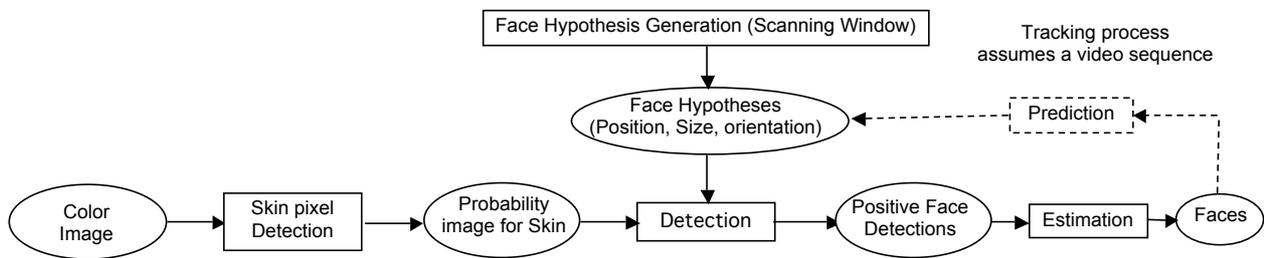
Covariance. Specifies the size and orientation of the mask.

1. Lab 1: Face Detection using Color

Skin color can be used to construct a simple detector for skin pixels in images. Color skin pixels can then be used to detect and track “blobs” that represent faces, hands and other skin colored regions in images.

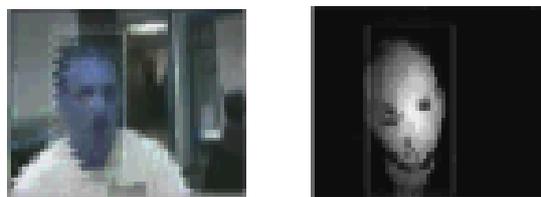
1.1 Algorithm Overview

The detector works by first computing the probability that each pixel contains skin. A sliding window (Region of Interest or ROI) is then scanned over the image. At each position, a weighted sum of probabilities is determined. Regions for which the weighted sum is above threshold are detected as faces. Adjacent face detections are grouped to form a single face detection.



Algorithm:

- 1) Compute probability of skin at each pixel.
- 2) Test for faces at possible positions and sizes (scanning Window).
- 3) Cluster adjacent detections.
- 4) Estimate precise face position, size and orientation from clusters of detections.



(images from a Robust Bayesian skin tracking in real time – Karl Schwerdt - 1993)

This system used detection with a Gaussian Window for robust tracking as described below.

2. Challenge 1: Detecting skin pixels with color

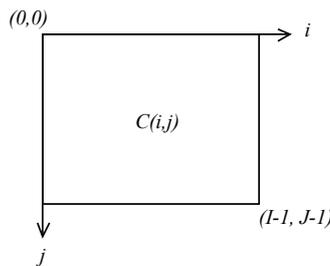
2.1 Baseline: RGB histograms

In the following, assume that we have a color image, where each pixel (i,j) is a color vector, $C(i,j)$, composed of 3 integers between 0 and 255 representing Red, Green and Blue.

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

Each color is represented by value between 0 and 255 (8 bits).

Let i refer to columns from 0 to $I-1$, and j refer to rows from 0 to $J-1$. Typical values for I and J would be 1024.



Suppose that we have K color images of size $I \times J$ pixels, $C_k(i,j)$.

This gives a total of $M = I \times J \times K$ pixels. Call this set of pixels X_m . (*Note: Do not confuse the image $C_k(i,j)$ with the class label C_k used in the last lecture.*)

Suppose that we have a ground truth function $y(X_m)$ that tells us whether each pixel is target (P or Postive) or not target (N or Negative). For FDDDB, this can be simply pixel inside a face ellipse as defined in the test data.

Suppose that a subset of M_T pixels that belong to a target class, T .

We allocate can two tables $h(r, g, b)$ and $h_T(r, g, b)$ and use these to construct two histograms.

$$\forall_{i,j,k} h(C_k(i,j)) = h(C_k(i,j)) + 1 ; M \leftarrow M + 1$$

$$\forall_{(i,j,k)} y(C_k(i,j)) = P : h_T(C_k(i,j)) = h_T(C_k(i,j)) + 1 ; M_T \leftarrow M_T + 1$$

For a color vector, $\vec{c} = \begin{pmatrix} R \\ G \\ B \end{pmatrix}$ we 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:

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

$P(T)$ is the probability that any pixel belongs to the target class. This can be estimated from the training data 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.

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

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

if $h(\vec{c}) = 0$ then $h_T(\vec{c}) = 0$ because $h_T(\vec{c})$ is a subset of $h(\vec{c})$. (we will need to test this case).

If we ASSUME that a new image, $C(i,j)$, 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 $P(i,j)$ that the pixel (i,j) belongs to the target T.

$$P(i,j) = L_T(C(i,j))$$

In this example, $h(\vec{c})$ and $h_T(\vec{c})$ are composed of $2^8 \cdot 2^8 \cdot 2^8 = 2^{24}$ cells. We define this as the Capacity of the histogram, Q.

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

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

This can result in a very sparse histogram. The reliability can be improved by using more training images from more cases.

A naive statistics view says to have at least 10 training samples for histogram cell.

That is $M \geq 10 Q$. However, it is often easier to work with powers of 2.

For example, $2^3 = 8 \approx 10$ which is approximately 10.

This suggest that we required $M \geq 8Q = 2^3 Q$.

Thus we would need $2^3 \cdot 2^{24} = 2^{27}$ training pixels. 2⁷ Meg.

(Note that a 1024 x 1024 image contains 2²⁰ pixels. This is the definition of 1 Meg)

Obtaining 2⁷ Meg pixels is not a problem for $P(\bar{c}) = \frac{1}{M} h(\bar{c})$ but may be a problem for training the histogram of target pixels $P(\bar{c} | T) = \frac{1}{M_T} h_T(\bar{c})$.

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? Two approaches:

We can reduce the number of values, N , for each feature, or we can reduce the number of features.

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".

2.2 Variation: Detection with Chrominance (normalized color).

Luminance captures local surface orientation (3D shape) while Chrominance is a signature for object pigment (identity). The color of pigment for any individual is generally constant. Luminance can change with pigment density (eg. Lips), and skin surface orientation, but chrominance will remain invariant.

Several methods exist to transform the (RGB) color pixels into a color space that separates Luminance from Chrominance.

$$\begin{pmatrix} L \\ c_1 \\ c_2 \end{pmatrix} \Leftarrow \begin{pmatrix} R \\ G \\ B \end{pmatrix}$$

A popular space for skin detection is computed by dividing R and G by luminance. These are often called "r" and "g" in the literature.

Luminance: $L = R + G + B$

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

The terms r and g have values between 0 and 1. To count statistics we need integer values. For this, it is common to represent r and g as integer numbers coded with N values between 0 and N – 1 by :

$$r = \text{trunc}\left((N - 1) \cdot \frac{R}{R + G + B}\right) \quad g = \text{trunc}\left((N - 1) \cdot \frac{G}{R + G + B}\right)$$

From experience, N = 32 color values seems to work well for skin with most web cameras, but this may change for certain data sets and should be experimentally verified.

Thus we can use a normalized vector $\vec{C} = \begin{pmatrix} r \\ g \end{pmatrix}$ 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 $R \times C$ where each pixel is an RGB color vector. This gives a total of $M = K \times I \times J$ color pixels.

Suppose that M_T of these are labeled as skin pixels i.e. $y(C_k(i,j)) = P$

We allocate two tables: $h(r,g)$ and $h_T(r,g)$ of size $N \times N$.

As before:

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

BEGIN

$$r = \text{trunc}\left((N-1) \cdot \frac{R}{R+G+B}\right) \quad g = \text{trunc}\left((N-1) \cdot \frac{G}{R+G+B}\right)$$

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

$$\text{IF } (y(C_k(i,j)) == P) \text{ THEN } h_T(r,g) = h_T(r,g) + 1 ; M_T \leftarrow M_T + 1$$

END

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

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

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

$$r = \text{trunc}\left((N-1) \cdot \frac{R}{R+G+B}\right) \quad g = \text{trunc}\left((N-1) \cdot \frac{G}{R+G+B}\right)$$

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

Where $P(i,j)$ is an image in which each pixel is a probability value from 0 to 1.

Probabilities can be expressed, of course be expressed as integers by multiplying by a quantization value (N).

2.3 Free Parameters

A number of parameters remain unspecified. Specification of these parameters generally depends on the application domain and the training data. Often these parameters make it possible to trade off error rates for computing time. To properly evaluate the algorithm, you should measure performance and compare performance over a range of parameters.

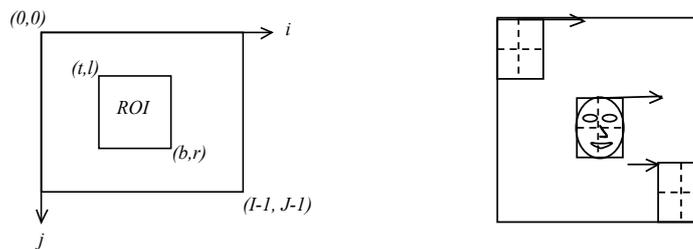
Free parameters for color based skin detection include:

- The use of color coding (eg, RGB, rg, other codings for chrominance)
- Value of N – the quantification for color color values.
- Training Data – different training and test regimes with variations such as the number of folds, and whether to train with a subset of the folds or all folds.

3. Challenge 2: Detecting Faces with Color

We can detect faces from the mass of skin detections within a “Region of Interests” (ROI). The ROI is determined either by a-prior knowledge or by some form of search procedure. In many cases this is based on tracking of targets in the scene. This is not possible for our example, because our training and test data are restricted to static images.

In many vision detectors the ROI is simply a sliding window that scans the entire image. This is the technique used, for example, with the Viola Jones Face Detector.



Generally both the size and the position of the ROI are scanned. Size can be on a linear scale or a logarithmic scale. There are good reasons to use a logarithmic scale for size (e.g. Fractional powers of 2), for covering large ranges of scales, but this is beyond the scope of this lecture.

A common representation for the ROI is a rectangle represented by four coordinates: (top, left, bottom, right) or (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 $(l,t), (r,b)$.

The center position, width and height are simply

$$c_i = \frac{l+r-1}{2}, \quad c_j = \frac{b+t-1}{2}, \quad w = l-r+1, \quad h = b-t+1$$

In many cases it is more convenient to fix c_i, c_j, w, h and calculate (l, t, r, b)

In this case the bounding box ROI is:

$$t = c_j - \frac{h-1}{2}, \quad b = c_j + \frac{h-1}{2}, \quad l = c_i - \frac{w-1}{2}, \quad r = c_i + \frac{w-1}{2}$$

The ratio of w and h can be fixed or variable. It is convenient to assure that w and h are odd numbers so that the boundaries fall on an integer pixel.

For detection in static images, we will need to test a range of ROIs with different positions and sizes. Note that with sufficient computing power, all windows can be processed in parallel.

In a scanning window system, the scanning step size, s , may also be varied. This may be the same size for row and column direction, or different step sizes may be used (say s_i and s_j)

Many engineers assume a step size of 1 pixel. However when the discriminant function is highly correlated at adjacent pixels (as is the case with color skin detection), then the resulting detection is very smooth. In this case a larger step size may be sufficient.

The largest reasonable value for s is half the width (and half the height if different step sizes in row and column are used). Beyond this the detection degrades rapidly because of aliasing effects, as can be demonstrated using signal processing techniques that are beyond the scope of this class.

When computing time is an issue, the detection algorithm can be optimized using a hierarchical search, starting with a large step size (Typically a power of 2) and recursively searching at higher step sizes near detected points.

3.2 Baseline: Sliding Window Detector

We can compute the ROI parameters for a face as a bounding box for an ellipse.

Let us define a face hypothesis as $\vec{X} = \begin{pmatrix} c_i \\ c_j \\ w \\ h \end{pmatrix}$

Since faces are generally vertical, we may assume that the ellipse is oriented with the major axis aligned with the row direction (j). The bounding box ROI is:

$$t = c_j - \frac{h-1}{2}, \quad b = c_j + \frac{h-1}{2}, \quad l = c_i - \frac{w-1}{2}, \quad r = c_i + \frac{w-1}{2}$$

The likelihood of a face at a position (c_i, c_j) of size (w, h) is:

$$g(\vec{X}) = \frac{1}{w \cdot h} \sum_{i=l}^r \sum_{j=t}^b P(i, j)$$

We can bias this likelihood to make a decision:

$$\text{IF } g(\vec{X}_m) + B > 0.5 \text{ THEN } R(\vec{X}_m) = P \text{ else } R(\vec{X}_m) = N$$

And of course

$$\text{IF } R(\vec{X}_m) = y(\vec{X}_m) \text{ THEN T else F.}$$

This technique will detect “faces” for a range of positions and sizes. As long as the detections overlap with the face ellipse in the data-base they are TRUE detections.

The problem with this technique is that faces are not square. The ROI gives equal weight to corners as the center. We can do better by weighting the pixels using a Gaussian function. This is called a “robust estimator” because it tends to reject outliers.

3.3 Variation: Robust Detection using a Gaussian mask.

In machine vision, fixation serves to reduce computational load and reduce errors by focusing processing on parts of the image that are most likely to contain information. A commonly practice is to use a Gaussian Window to suppress signals outside of the fixated region. A similar technique is used to suppress outliers for robust estimation.

A Gaussian window has the form:

$$G(i, j; \vec{\mu}, \Sigma) = \frac{1}{2\pi \det(\Sigma)^{1/2}} e^{-\frac{1}{2}(\vec{p}-\vec{\mu})^T \Sigma^{-1}(\vec{p}-\vec{\mu})}$$

Where $\vec{p} = \begin{pmatrix} i \\ j \end{pmatrix}$ represents image positions, and $\vec{\mu} = \begin{pmatrix} \mu_i \\ \mu_j \end{pmatrix}$ is the center position of the window.

The covariance matrix of the window is: $\Sigma = \begin{pmatrix} \sigma_i^2 & \sigma_{ij} \\ \sigma_{ij} & \sigma_j^2 \end{pmatrix}$

And $\det(\Sigma)$ is the determinant of the covariance, Σ . $\det(\Sigma) = \sigma_i^2 \sigma_j^2 - \sigma_{ij}^2$

The term $\frac{1}{2\pi \det(\Sigma)^{1/2}}$ assures that the window weights sum to 1.

We can use the parameters of the window to define a face hypothesis. $\vec{X} = \begin{pmatrix} \mu_i \\ \mu_j \\ \sigma_i^2 \\ \sigma_j^2 \\ \sigma_{ij} \end{pmatrix}$

Normally the Gaussian mask should be computed within a ROI determined by a bounding box. Typically the bounding box should be at least 3σ (standard deviations), but if speed is more important the false negatives, then computation time can be reduced by using a smaller ROI, down to as small as 1 standard deviation.

We can define the ROI as a 3σ bounding box:

$$t = c_j - 3\sigma_j, \quad b = c_j + 3\sigma_j, \quad l = c_i - 3\sigma_i, \quad r = c_i + 3\sigma_i$$

To evaluate a face hypothesis, we compute the face likelihood as a weighed sum of the skin probabilities by the Gaussian mask.

$$g(\vec{X}) = \sum_{i=1}^r \sum_{j=1}^b P(i, j) G(i, j; \vec{X})$$

As before, the discriminant, $g(\vec{X})$, has a value between 0 and 1. This is the likelihood that a face may be found at \vec{X} .

As before, faces are detected as:

$$\text{IF } g(\vec{X}_m) + B > 0.5 \text{ THEN } R(\vec{X}_m) = P \text{ else } R(\vec{X}_m) = N$$

3.4 Free parameters to test

As with color skin detection a number of parameters remain unspecified. To properly evaluate the algorithm, you should measure performance and compare performance over a range of parameters.

Free parameters for color face detection include:

For a simple scanning window, these include:

- Width and height of ROI
- Range of positions
- Step size for scanning windows
- The percentage of overlap with the ground truth that is considered a TRUE detection.

For a Gaussian detection window, parameters also include

- The width and height of the ROI compared to the standard deviations of the principal axis.
- The range and step sizes for orientation of the Gaussian window.

4. Challenge 3: Face Localization

4.1 Detection vs Localization (Parameter Estimation)

Detection: A face is considered to be “detected” if a positive detection is found at a position that overlaps the face in the ground truth. The simplest form of test is to verify that the location of the face is within the ellipse of the ground truth label for a face in the image.

For a face hypothesis at $\vec{X} = \begin{pmatrix} c_i \\ c_j \\ w \\ h \end{pmatrix}$ a face is detected if $R(\vec{X})=P$

The detection is TRUE if $R(\vec{X})=y(\vec{X})$

Using the ground truth $y(\vec{X}) = \begin{cases} P & \text{if } \left(\frac{(x-c_x)^2}{r_a^2} + \frac{(y-c_y)^2}{r_b^2} \leq 1 \right) \\ N & \text{otherwise} \end{cases}$

The Hypothesis is a TRUE POSITIVE detection if $y(\vec{X})=P$ and $R(\vec{X})=y(\vec{X})$.

A large number of TP faces will be detected for each ground truth face. If the search space includes multiple scales and orientations, than the number of detected faces will be even larger. All of these correspond to the same face!

Localization (more precisely parameter estimation) specifies precisely where, and with what parameters, the face may be found. There should be only ONE face located for each true face. The face should be located at a position, size and orientation as close to the true face as possible. When searching at multiple scales and orientations, each detection has the form of an error vector.

A distance metric, relating degrees and scale change to position is required to reduce this to a single number. For discrete samples the distance metric can provided implicitly by the sample step size in scale, orientation and position.

We can obtain a location (or parameter estimation) from multiple positive detections by suppressing detections for which the discriminant, $g(\vec{X})$ is not a local maximum. This requires specifying a measure for locality.

We can also estimate the parameters of the face from the moments of a “cloud” of detections at multiple positions, scales and orientations. This is the same robust estimation technique that we used above, extended to robustly estimate position, size and orientation.

4.2 Baseline: Localization by non-maximum suppression

In order to suppress non-maximal detections, the easiest method is to build a list of detections hypotheses, $\{\vec{X}\}$ over the desired range of positions, scales, orientations and any other parameters, and then filter this list to remove any hypothesis for which the discriminant not a local maximum. (A maximum within a some distance R.)

$\forall \vec{X}$

$\forall \vec{X}_i, \vec{X}_j \in \{\vec{X}\} : \text{IF } \text{Dist}(\vec{X}_i, \vec{X}_j) < R \text{ AND } g(\vec{X}_i) < g(\vec{X}_j) \text{ THEN } \{\vec{X}\} \leftarrow \{\vec{X}\} - \vec{X}_i$

This requires defining some notion of distance that includes position, scale and orientation. This is generally done with regard to an expected range of positions, orientations and scales over which a single face would be detected. For example, using a Mahalanobis distance (Distance normalized by covariance).

4.3 Variation: Localization by Robust Estimation

We can use robust estimation to estimate the most likely parameters from a set of detections. To do this, we will calculate the weighted moments from the set of detections.

Suppose that we have a set of N detections: $\{\vec{X}_n\}$ and that for each detection we have a discriminant $g(\vec{X}_n)$.

The mass of the detections is $M = \sum_{n=1}^N g(\vec{X}_n)$. This is the zeroth moment of the set of detections.

The "expected value" is $E\{\vec{X}\} = \frac{1}{M} \sum_{n=1}^N g(\vec{X}_n) \cdot \vec{X}_n$

This is the first moment (or center of gravity) of the values of $\{\vec{X}_n\}$.

The expected value is a vector of first moments for each parameter:

For D parameters, the center of gravity is a vector

$$\bar{\mu} = E\{\bar{X}\} = \frac{1}{M} \sum_{n=1}^N g(\bar{X}_n) \bar{X}_n = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_D \end{pmatrix} = \begin{pmatrix} \frac{1}{M} \sum_{n=1}^N g(\bar{X}_n) X_{1n} \\ \frac{1}{M} \sum_{n=1}^N g(\bar{X}_n) X_{2n} \\ \dots \\ \frac{1}{M} \sum_{n=1}^N g(\bar{X}_n) X_{Dn} \end{pmatrix}$$

The vector $\bar{\mu}$ the vector of weighted averages for the components of \bar{X}_n .

If there is only one face inside the set $\{\bar{X}_n\}$ of positive detections then $\bar{\mu}$ provides the most likely estimate for set of parameters the detection, (e.g. position, size and orientation).

In the case of multiple faces, the best estimate for each of the faces can be determined using the Expectation Maximization (EM) algorithm. Alternatively, the estimation may be limited to faces detections within a small region of the image.

4.4 Free parameters to test

As with color face detection a number of parameters remain unspecified. To properly evaluate the algorithm, you should measure performance and compare performance over a range of parameters.

Free parameters for color face detection include:

For a simple scanning window, these include:

- Number of Faces in the image
- Range of positions, size and orientations to test
- Distance to use for non-maximal suppression.
- Size of the region used to cluster faces detections.