

Bogdan KWOLEK

Rzeszów University of Technology

W. Pola 2, 35-959 Rzeszów, Poland

bkwolek@prz.rzeszow.pl

POINTING ARM POSTURE RECOGNIZING USING STEREO VISION SYSTEM

Abstract.

This paper describes an approach to real-time detecting static pointing arm-postures which can be used to send position instructions to a mobile robot. This takes place by analyzing the color image and utilizing the range image of the scene obtained by a stereo vision system. First the algorithm locates the person's face and afterwards the arm-postures are recognized. The start as well as the end of the pointing command have been signaled by suitable static arm-postures of an instructor who stands in front of the robot. The system has been tested in tasks consisting in selecting the targets on the laboratory floor. All image processing is realized on-board the robot.

keywords : *Color Image Processing, Vision-Based User Interfaces, Robot Vision.*

1 INTRODUCTION

Arm-postures can be a very useful and compact way to convey geometrical information to the robot [8]. One of the most expressively significant arm-postures is pointing towards an object [6], [16]. Having a robot with a recognition interface, the user can point to the desired spot on the floor, rather than input the location to where the robot should move. For non-expert users it is easier to point to an object than employ the verbal description of its coordinates.

The aim of this work was to build a system that makes possible realization of experiments which consist in determining which object from among those lying on the floor and being in the view of the camera the user has pointed to, giving commands in front of the autonomous mobile robot. The recognition of such commands has been realized on the basis of template matching techniques [15]. The centers of correlated sub-images with the manually prepared in advance templates have been located in the area which location had been determined on the

basis of eyes position. The determination of eyes position has in term been preceded by a searching within skin-like regions for symmetrical face candidates as well as candidates of face which match suitable eyes-template. The final verification of face presence has been realized using the well known eigenfaces approach [18]. The information about the relative distance of the user's face to the camera has allowed us to avoid search over scale parameter in the matching procedure used for arm-posture recognition. The determination of floor spot which was pointed to by user and of the mentioned above scale ratio as well as the selection of the eyes-template size have been realized on the basis of stereo vision.

The contribution of this work relies on the employment of several visual cues and in particular stereo with the aim to detect and to localize the face and then to recognize in real-time static pointing arm-postures. By using the stereo system the positions of the person's face and hand in the scene are determined and used then by module responsible for determination of the direction of pointing. Since false-positives (the robot recognizes the command which was not given by the user) are in human-machine communication less desirable than false-negatives (the robot fails to recognize the hand-action), our approach rests on reliable face detection and localization.

Human face detection plays an important role in human interface. Some of the best techniques for face detection are still too computationally expensive and cannot match the real-time constraint using the personal computers that we have today. The eigenface approach is one of the few techniques that can fulfill the real-time requirements. However, most eigenface-based systems work with the assumption that the location of a face within image is known [23]. The use of pattern matching with the patterns being prepared manually is reported in work [2]. Since the face size and position were not known a priori the searching has been realized on all possible sizes and locations of the head. The system reported in [2] takes 2 seconds on a SGI Indigo 2. A system for real-time face detection and tracking on the basis of pan-tilt-zoom controllable camera has been presented in [21]. The face detection takes about 0.2 seconds on a Pentium 266 MHz. Detected skin pixels are grouped into regions and then regarded as faces if dark regions of eyes, eyebrows and mouth are detected within these regions. Sobotkka at al. [14] showed that human skin color for all races are clustered in normalized RGB space. In work [12] shape symmetry is used to verify whether it is correct face region or not. In work [20] a robot looks for dark facial regions within face candidates obtained on the basis of color. Distance to the tracked color is roughly determined by measuring the area covered by the color in the image. A binary matching technique is used to detect the face in area located above of the rectangular boundary of the detected color shirt. The sizes of sub-images are scaled according to the robot to person distance provided by the sonar sensors. In work [19] a fast and adaptive algorithm for tracking and following a person by robot-mounted camera has been described. After locating the person in front of the robot an initial probabilistic models of the person's shirt and face colors are created. The system uses window of fixed size to adapt the face and shirt color models over time and it is essential that these windows can only contain face and shirt colors. In particularly, the distance between the robot and the person must be approximately fixed. Two alternative methods for pose analysis are used: neural networks and template matching. The Viterbi algorithm is utilized to recognize motion gestures. The Perseus system [6] is capable of finding the object pointed to by a person. The system assumes that people is only moving object in the scene. Perseus uses independent types of information e.g. feature maps and disparity feature maps. The distance between the robot and person may vary. In other system [13] a behavior for a

person following with an active camera mounted on robot has been presented. In that system the head of person is located using skin color detection in the HSV color space and color thresholding techniques.

The organization of the paper is as follows. Color image processing is discussed in section 2. In section 3 the algorithm for face presence detection and face localization is described. Arm-posture recognition is presented in section 4. In section 5. experimental results are briefly overviewed. Finally, in section 6 a conclusion follows.

2 COLOR IMAGE PROCESSING

Color image segmentation can be useful in many applications. On the basis of segmentation results, it is possible to identify regions of interest and desirable objects. The problem of segmentation is not easy because of image texture, shadows, etc. If an image contains only homogeneous color regions, clustering methods are sufficient to handle the problem. In practice, natural scenes are rich both in color and texture as well as they are under varying lighting conditions. Therefore our approach uses color only for coarse extraction of regions of interest and additionally utilizes the range information. In particular, the range information allows us to reduce search over scale or select appropriate template size in correlation procedure. Detection on the basis of matching is an expensive operation, involving search over the input image and over the scale during matching with the model. To reduce the working area in the matching procedure utilized in arm-posture recognition we have chosen the eyes location as the datum point. The search for eyes has been preceded by skin regions detection.

Face color as a feature for tracking people has been used for example in work [22]. The advantage of such approach relies on the speed at which present low-cost personal computers can extract the object of interest from color images. That aspect is particularly important considering on the one hand the limited computational power of on-board computer of the mobile robot and on the other hand the necessity of work at rate of around 10 Hz.

The most common space used in computer displays and CCD cameras is the RGB color space. Since the brightness can change rapidly when robot moves, most of the known color segmentation algorithms operate using chromatic colors. The authors of work [1] came to the conclusion that the best space for skin color representation is the normalized color space. The two-dimensional intensity-normalized chromaticity color space is relatively robust to the change of the illumination. But an obvious shortcoming is that the pure colors are unstable and meaningless when $R+G+B$ is small [10]. Therefore following a suggestion from work [17] we have only employed the chromatic areas of image in further color processing.

The skin color distribution in normalized space of colors can be represented by a Gaussian model [22] $M = (\mu_r, \mu_g, \Sigma)$, where μ_r, μ_g are the means and Σ is the covariance matrix. The parameters of Gaussian model have been obtained on the basis a set training images. Each image with user outlined regions of interest has undergone the color space transformation. The aim of such an operation was to obtain a representative sample of typical skin colors of objects and thus it was not essential to outline the whole face. One attraction of bimodal normal distribution is that it can be used to generalize about small amounts of training data. To extract the face candidates, each pixel was examined by Gaussian model.

The better the pixel matches the color model, the higher the probability and response of such a color filter is, see fig. 1b.

Filtering on the basis of the mask whose weights correspond to the binomial coefficients [4] was the first operation realized on acquired images. Then the filtered RGB values were transformed into probability ones. To shorten the execution time the lookup table for direct conversion from RGB to probability was utilized.

A pixel was identified as a candidate of the face if the corresponding probability was above a threshold. The threshold has been set up to a low value. Because we have been interested in coarse extraction of face candidates, the choice of threshold did not have a considerable influence on obtained results. After an image was thresholded, morphological closing operator (dilate followed by erode) [5] was performed. The aim of such an operation was to fill small holes and to smooth border.

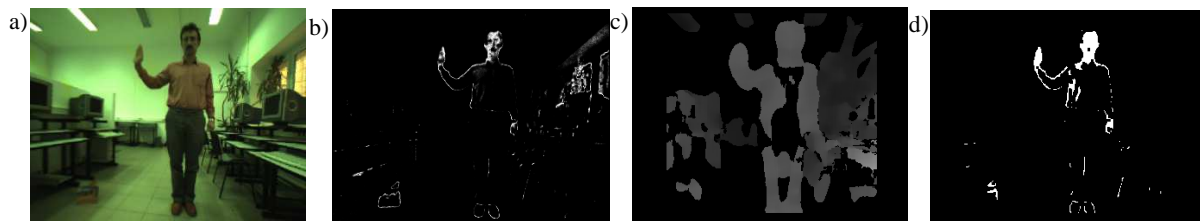


Fig. 1 Face candidates extraction. Raw color image (a). Probability image (b). Depth image (c). Closed image (d)

Our system utilizes the depth information provided by the commercial SRI Megapixel Stereo Head [7]. Corresponding pixels between images are detected using the area-correlation method where the correlation is used to obtain the most correspondences between small patches [7]. We have set up the number of disparities to 32 and thus worked with a minimal horopter distance of 96 cm and a maximum range resolution of 7.5 mm. Stereo depth information that was utilized at this stage can be found in small area correspondences between image pairs and therefore it gives poor results in regions of little texture, see fig. 1c. However, the depth map covering a face region is usually dense because the human face is rich in details. The person performing hand-actions has stayed about 1.5 m in front of the robot. Taking the depth image, we have removed from closed image the pixels representing a distance to the camera greater than 2.5 m, see fig. 1d. The image prepared in such a way was then labeled. The objective of the realization of this operation was to determine distance of each individual scene object to the camera.

The connected component analysis was utilized to gather neighborhood pixels into larger regions. Connected component labeling scans an image, pixel-by-pixel, from top to bottom and from left to right in order to extract pixels and regions which share the appropriate set of intensity values. Each object receives a unique label which can be used to distinguish regions during further processing. In our approach the connected components algorithm works on a binary image and is based on region identification in 8-connectivity [15]. At each point the algorithm examines pixel p from binary image and if that pixel is set to foreground, it checks its four neighbors from labeled image which have already been processed (i.e. the neighbor to the left of actually analyzed pixel, above it, upper left and right). Having obtained values of pixels, one of the following can occur:

- if all four neighbors are 0, assign a new label to pixel p , else
- if only one neighbor has a label, assign its label to pixel p , else

if one or more of the neighbors have two different labels, assign to p a label which is smaller and update the glue table for equivalencies.

Each pixel is replaced in the second pass by one of the unique labels which had been assigned to each region and stored in the glue table in the previous pass. The labels obtained in such a way were then used within the depth image to determine an average range of each labeled region.

3 FACE PRESENCE DETECTION AND FACE LOCALIZATION

The matching procedure which was used for arm-posture recognition has been preceded by face localization module. The reliable face location has allowed us to decrease the number of sub-images with expected arm-postures because the operating area of matching procedure was bounded to the region that surely contained the arm and body of the user. Thanks to such a datum point we had to re-center the arm-templates in relative small surrounding and thus the recognition was realized fast and reliable. Stereo range gives information how distant objects are and thus for objects of known size such as people who take place in our experiments we could set the image scale for arm-postures matching.

Face detection methods can be divided into four groups [23]. In knowledge-based model techniques, the model consists of features and their relation between each other is derived from our knowledge of human faces. The relations of features are expressed as the sets of rules and the verification of hypotheses on the basis of a decision tree is integrated within a search for a face. The methods based on feature-invariant approach also search for facial features. In this approach the face candidates are extracted by maximizing search criteria. The set of face features is assumed to be invariant regardless of the face viewpoint, orientation. These features can be geometric or skin color (texture) based. The approach which is based on template matching maximizes a correlation function of a human face template over the whole image. Sample image window with a face appearance is sometimes normalized and scaled. The template is correlated and in the point where the correlation has the largest value the hypothesis is then formulated. The appearance-based methods operate in a multidimensional feature space. In the methods that are based on eigenfaces the feature space is reduced using principal components. The scaled query windows of the input image are projected into classification subspace and the closer the distance to the projected training image is, the greater the probability that such window corresponds to the trained human face is. Our approach to face detection relies on mixture of the last three techniques. After face detection we know exact three-dimensional position of face with respect to the camera.

The first stage of our algorithm for face localization operates in gray images and it searches for symmetrical face candidates. The vertical axis of symmetry has only been considered. Thanks to the being in disposal information about distance of the examined candidate of face to the camera, the reflected symmetry has been checked in windows with dimensions covering only the face. The non-zero pixels of closed image have only been considered as centers of windows. Next, taking again into consideration the range information, the eyes-template of appropriate size has only been applied in the centers of symmetry detected at previous stage. Templates that model human eyes have been formed based on real face images of different persons. The arrangement of elements in sample

template is shown in fig. 2. The gray values in applied template symbolize the regions not considered in matching process, whereas the white and black areas symbolize regions which

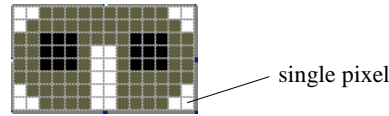


Fig. 2 Matching template used in searching for eyes candidates

should be above or under the mean value, respectively. The template was compared, pixel by pixel, to pixel blocks recentered in the centers of symmetry. A pixel block was considered as eyes if the number of matching pixels was greater than a threshold. The locations detected in such a way were then used by the face presence detection procedure that was based on the eigenfaces algorithm.

The aim of the usage of the eigenfaces was to reduce false-positive errors of face localization. The arm-postures have been recognized by matching procedure on binary images. In the case of false-positive localization of face, the localization of matching template is incorrect, and thus the output value of matching process can reflect commands which were not given by the user. Therefore the reliable face detection and localization play an important role in our algorithm of given hand-postures recognition. The reliable face detection can be very useful in several robotic tasks [3].

The eigenfaces algorithm operates on gray images and collection of training faces. We have prepared a collection of training faces that consists of 64 images. While preprocessing, the average image for this collection was computed. Next, this image was subtracted from each training image and placed in the matrix. This matrix was used to create the covariance matrix. The eigenvalues and eigenvectors of such a covariance matrix were determined. The first 16 normalized eigenvectors, sorted by decreasing eigenvalue represent subspace in which classification at run-time phase was performed. The eigenfaces are normalized eigenvectors which are the principal components of face space and they reflect the statistical properties of facial appearance. The first ten eigenfaces related to our training collection are shown in fig. 3.



Fig. 3 The first ten eigenfaces computed from our training collection

Prior to the run-time phase the set of reference images was read and projected into the classification subspace. In run-time the L1 distance between these images and the projected (and centered) image onto the subspace was computed to determine the closest match. The size of the sub-image which has to be cropped was scaled according to the distance from the robot to the person. The location of the sub-image was determined on the basis of eyes location.

The eigenfaces method is more robust than other template matching techniques that are based on detection of visible local features as well as distances between them. The algorithms based on eigenfaces recognize only a single face appearance that has been taken from a narrow angle, most often in the frontal view. Thanks to the obtained in advance localization of face candidates, the presence of a vertical and frontal-view face in the scene was verified in a very fast manner.

4 ARM-POSTURE RECOGNITION

The arm-postures that have been used by our system can be divided into two categories, i.e. signaling the start as well as the end of the pointing command and the proper pointing command. The recognition of the first category command relies on the matching templates, see fig. 4. b, c. The area which was not considered in the matching process was marked with suitable gray value. The body of the user and its surrounding background were marked by white and black values, respectively. The arm and its surrounding background received gray values that are close to white and black values, respectively, see fig. 4. b, c. The five weights inversely proportional to the occupied areas were used during the matching process with the mentioned templates.

The region growing is the process of grouping adjacent pixels or a collection of pixels which share similar attributes into larger regions [5]. The algorithm starts with a number of seed pixels and then from those it expands the regions by adding unassigned adjacent pixels that satisfy a desired criterion of similarity with the seed pixels. This algorithm was applied in our approach to extract the shirt of the user. The seed region was located in the area below the face where a depth compared with the depth of the face was detected. The following homogeneity criteria were used in a two-stage region growing

$$\sqrt{(I - \bar{I})^2 + S^2 + \bar{S}^2 - 2S\bar{S}\cos(H - \bar{H})} \leq d_1, \quad |D - \bar{D}| \leq d_2 \quad (1)$$

where H, S, I (hue, saturation, intensity) are current values of color components of tested pixel, D is depth, $\bar{H}, \bar{S}, \bar{I}, \bar{D}$ are mean values related to the seed region, d_1 and d_2 are threshold values. The HSI criterion considers all three color components and takes into account the cylindrical nature of HS components. For reasons of angular value of hue the values belonging to the seed region have been converted into Cartesian coordinates in the following manner:

$$x_i = \cos(H), \quad y_i = \sin(H) \quad (2)$$

The average values $\bar{x} = \frac{1}{N} \sum_{i=0}^{N-1} x_i$ and $\bar{y} = \frac{1}{N} \sum_{i=0}^{N-1} y_i$ of the seed region were used in the arctan function for computing the mean \bar{H} value. After the first stage of the region growing was finished, the average values $\bar{H}, \bar{S}, \bar{I}$ of the detected region were extracted and then in the second stage the HSI criterion with obtained in this manner new seed region and reduced to half d_1 value were used. The binary image which corresponds to the extracted regions of interest is presented in fig. 4a. The binary sub-image containing the user's body was scaled before matching according to the robot to the person distance. The only disadvantage of the HSI color space is the costly conversion from RGB color space. We handled this problem by using lookup tables.

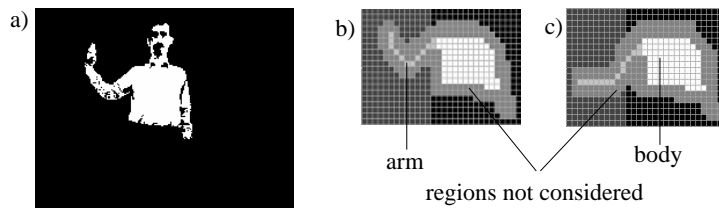


Fig. 4 Sample binary image used in arm-posture recognition (a). Templates (b), (c)

Once the start command had been executed, the pointing arm-posture was watched for. During the pointing phase the arm remains stationary for a short time. Once the face position was known, the procedure responsible for determination of direction of pointing looked for the person's hand, which were assumed to be located to the right of the face. As with the face detection, the position of the hand was determined on the basis of labeled as well as depth images. Next, the consistency checks have been realized to verify that face-hand relationship which can be characterized by distances, angles, etc., is physically possible. The distances of the face and the hand to the camera that had been obtained from the depth image as well as positions of them allowed us to obtain direction of pointing. Then a line from the person's face center to the pointing hand was found and a cone was centered around it, starting at the hand. We have assumed that between the start command and the pointing command the user occupies approximately the same position. Taking the above into consideration, we obtained the face center on the basis of the skin-color region which was situated near the detected one during the execution of the start command, see fig. 5.

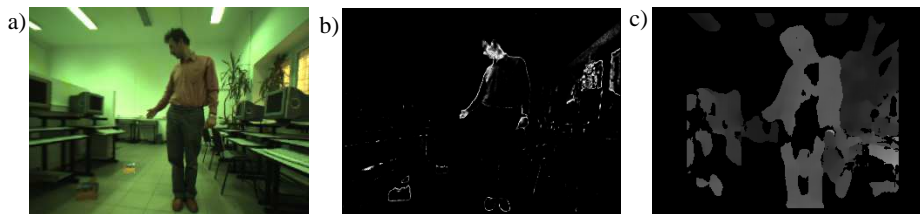


Fig. 5 Pointing arm-posture recognition. The raw input image (a). Probability image (b). Depth image (c)

5 EXPERIMENTAL RESULTS

We have presented an algorithm which was implemented in system [9] and tested in several experiments. We have tested the algorithm on three people. The spatial resolution of pointing is about 20 cm and when the targets are at least 20 cm apart we can correctly tell a two targets apart in over 90% cases. The algorithm runs at 320x240 image resolution at frame rates of 6-10 Hz on 850 MHz Pentium III laptop which was installed on Pioneer 2 DX mobile robot [11].

6 CONCLUSION

We have presented a method to recognize the pointing arm-posture which can be used in tasks including robot programming by demonstration. The recognition has been performed on the basis of the information obtained from stereo cameras and the analysis of color as well as gray images. The combination of low-level image processing algorithms and of generic image features has allowed us to recognize the pointing command fast and robustly. One of the individual particularities of the presented method that makes it reliable is that it utilizes the information about the presence of the user's face. The depth map covering the face region is usually dense and this together with skin-color and symmetry information as well as eyes-

template assorted with the depth has allowed us to apply the eigenfaces method and to detect the presence of the vertical and frontal-view faces in the scene very reliably.

ACKNOWLEDGEMENTS

This work has been supported by the Polish Committee for Scientific Research (KBN) within the project 7T11C03420

REFERENCES

- [1] L. M. Bergasa, M. Mazo, A. Gardel, M. A. Sotelo, L. Boquete, *Unsupervised and adaptive Gaussian skin-color model*, Image and Vision Computing, 18, 2000, 987-1003.
- [2] Q. Chen, H. Wu, M. Yachida, *Face detection by fuzzy pattern matching*, In Proc. 5th IEEE Int. Conf. Comput. Vision, Boston, 1995, 591-596.
- [3] S. Feyrer, A. Zell, *Detection, tracking, and pursuit of humans with an autonomous mobile robot*, In Proc. Int. Conf. on Intell. Robots and Syst., Kyongju Korea, 1999, 864-869.
- [4] J. M. Gauch, *Noise removal and contrast enhancement*, In: S. J. Sangwine, R. E. N. Horne (eds.): The Colour Image Processing Handbook, Chapman & Hall, London, 1998.
- [5] B. Jähne, *Digitale Bildverarbeitung*, Springer-Verlag, Berlin Heidelberg, 1997.
- [6] R. E. Kahn, M. J. Swain, P. N. Prokopowicz, R. J. Firby, *Gesture recognition using the Perseus architecture*, In Proc. of the IEEE Conf. on Computer Vision and Pattern Rec., San Francisco CA, 1996, 734-741.
- [7] K. Konolige, *Small Vision System: Hardware and implementation*, In Proc. Int. Symposium on Robotics Research, Hayama Japan, 1997, 111-116.
- [8] D. Kortenkamp, E. Huber, R. P. Bonasso, *Recognizing gestures on a mobile robot*, In Proc. of AAAI, 1996, 915-921.
- [9] B. Kwolek, *Person following, obstacle detection and collision-free path planning on autonomous mobile robot with color monocular vision*, In Proc. of the 8th IEEE Int. Conf. Methods and Models in Automation and Robotics, Szczecin, 2002, vol. II, 971-978.
- [10] Y. I. Ohta, T. Kanade, T. Sakai, *Color information for region segmentation*, Computer Graphics and Image Processing, 13, 1980, 222-241.
- [11] *Pioneer 2 mobile robots - Operations Manual*, ActivMedia Robotics, LLC, 2001.
- [12] E. Saber, A. M. Tekalp, *Frontal-view face detection and facial feature extraction using color, shape symmetry based cost functions*, Pattern Recognition Letters, vol. 19, 1998, 669-680.

-
- [13] H. Sidenbladh, D. Kragić, H. I. Christensen, *A person following behaviour for a mobile robot*, In Proc. of the IEEE Int. Conf. on Robotics and Automation, Detroit MI, 1999, 670-675.
- [14] K. Sobottka, I. Pitas, *Extraction of facial regions and features using colour and shape information*, In Proc. of ICIP, 1996, vol. III, 483-486.
- [15] M. Sonka, V. Hlavac, *Image processing, analysis and machine vision*, Chapman & Hall Comp., London, 1994.
- [16] J. Triesch, Ch. von der Malsburg, *A gesture interface for human-robot-interaction*, In Proc. the IEEE 3rd Int. Conf. on Aut. Face and Gesture Recognition, Nara Japan, 1998, 14-16.
- [17] D. C. Tseng, C. H. Chang, *Color segmentation using perceptual attributes*, In Proc. 11th Int. Conf. on Pattern Rec., Den Hague Netherlands, 1992, vol. III, 228-231.
- [18] M. A. Turk, A. P. Pentland, *Face recognition using eigenfaces*, In Proc. of Conf. on Computer Vision and Pattern Recognition, 1991, 586-591.
- [19] S. Waldherr, R. Romero, S. Thrun, *A gesture based interface for human-robot interaction*, Autonomous Robots, vol. 9, 2000, 151-173.
- [20] C. Wong, D. Kortenkamp, M. Speich, *A mobile robot that recognizes people*, In Proc. 7th IEEE Int. Conf. on Tools with Artificial Intelligence, 1995.
- [21] G. Xu, T. Sugimoto, *A software-based system for realtime face detection and tracking using pan-tilt-zoom controllable camera*, In Proc. of Fourteenth Int. Conf. on Pattern Recognition, 1998, vol. II, 1194-1197.
- [22] J. Yang, A. Waibel, *A real-time face tracker*, In 3rd IEEE Workshop on Appl. of Computer Vision, Sarasota Florida, 1996, 142-147.
- [23] M. H. Yang, D. Kriegman, N. Ahuja, *Detecting faces in images: A survey*, IEEE Trans. on Pattern Analysis and Machine Intelligence PAMI, vol. 24, no. 1, 2002, 34-58.