

General Framework for Multi-Player Tracking In Volleyball Sports Videos

Hananeh Salehifar

Faculty of Electrical, Computer, and IT
Islamic Azad University, Qazvin Branch,
Qazvin, Iran
H.Salehifar@qiau.ac.ir

Azam Bastanfard

Dep. of Computer, Faculty of Engineering
Islamic Azad University, Karaj Branch,
Karaj, Iran
Bastanfard@kiau.ac.ir

Abstract— The task of reliable detection and tracking of multiple players becomes highly complex for indoor team sports. In this paper, a novel robust framework is presented for multi-player tracking, using a static camera in volleyball sports video. Proposed framework is based on combination of GMM, Mean Shift and Kalman Filter. Kalman Filter is used to predict the next position of every player in each frame. The place of player is calculated by Mean Shift and GMM algorithms and then the existing noise will be removed by Kalman Filter. In comparison to several similar trackers the proposed tracker produced better estimates of position and prediction as well as reducing the number of failures.

Keywords- *player tracking; Kalman filter; volleyball videos;*

I. INTRODUCTION

Video event tracking for understanding, retrieval and search is becoming more and more important recently. One of the interesting applications is intelligent sports video summarization and browsing based on semantic analysis, e.g., segmenting video sequences into camera shots, tracking players, recognizing players' intentions and analyzing the overall play process, mapping low-level features to high-level events (i.e., classification and indexing) and generating natural language description. For analyzing the players' activity in different sports, information is needed about movement of players participating in matches. Players' motion data reveals many aspects that are not directly visible. For instance, it highlights the reasons why some players perform better than others and it suggests methods of training to make good players perform even better. Locating, labeling and tracking players have broad applications, especially for a team analyzer and in the broadcast sports videos industry. This task is quite challenging because of many difficulties, such as occlusion, similar player appearance with low discrimination, varying number of players, abrupt camera motion, player pose variance, various noises and video motion blur. The development and use of a video-based, computer-assisted track-and-field motion acquisition system is at the center of our research. It delivers the data about players' position, velocity and acceleration and conforms to the requirements of sports analysis. Many commercial sports analysis systems have been developed in recent years. But many of these systems require some kinds of markers to be attached to the body of the athlete, which are distracting and not acceptable during regular league or championship matches. Besides, many researchers paid their

attention to automatic annotation systems of football matches. For merely, researchers used self-made, video-based systems with a resolution of 1 m for successful analysis of a soccer match [1]. In 1996, Needham [2] proposed a soccer player tracking system over a large playing area from a single fixed camera, and he also developed his method into a multi-camera system [3]. In recent years, research on soccer video analysis [4–6] has been growing, and various effective methods have been applied. For instance, Intille and Bobick [7] studied large-scale motion acquisition and analysis, with the ultimate goal of realizing a fully automated tracking of American football players. Unfortunately, many important aspects of motion acquisition were neglected, and these researches are too domain specific to be applied in track and field sports directly. There are also computer vision systems developed in other sports domains. In [8–10] they used two stationary cameras mounted directly above the court and proposed a handball player's tracking approach. Similar methods to theirs were applied to other sports such as squash [11] and basketball [12]. Okuma [13] developed a hockey annotation system to automatically analyze hockey players. However, frequent occlusion has not been considered. Volleyball is of the plays, which is rarely studied. The excess of players in proportion to the court, rapid motion of players, wide range of possible motions change, non rigid shape, and frequent occurrence of occlusion among players are of the items which make it hard to analyze volleyball. Thomas and Christiana [14] designed a beach volleyball motion analysis method using a single camera. Every team has two players in beach volleyball. The presented system studies tracking two players in front of camera, and since there is little occlusion between two players, this system has not studied the most important problem of analyzing the plays and also the volleyball, which is occlusion. The number of beach volleyball players are not comparable with indoor volleyball, because every indoor volleyball has six players. So it enhances the complexity of analyze and also the presented system for indoor volleyball is not comprehensive. The problem of occlusion among players should be considered for this kind of play. Thomas and Peter have studied specific motion such as digging, running, over head passing and side away running in volleyball [15]. It is worthy to mention that these motions are not studied in real situation of volleyball and are performed by only two players. We have presented in our previous work [16] a method for tracking players based on adaptive background subtraction by two concurrent frames. In

this method, the difference between two successive frames is calculated and moving players are detectable. We have used a team playing in Iran volleyball league official match to test this method. Presented algorithm is an efficient method for tracking players, it also has some faults. One of the faults is the loss of recognition of some players who do not have any movements in some successive frames. The tracking will be disturbed if the player stops and doesn't move. Because this method is based on differences between similar pixels in successive frames. If sever occlusion happens between players, this method will be inefficient. In [17], we proposed a method about this subject, which is based on GMM and Kalman Filter. The mentioned method, separated players from background are calculated by GMM and next location of every players in next frames estimation via Kalman Filter. Although this method is better than previous one, the offered method is able to solve occlusion problem partially. Thus, to analyze the play, we need a comprehensive method that would have suitable efficiency in all complex cases. Therefore, presenting a general framework of volleyball in video images can help coaches, analysts and referees to study and recognize any kinds of fouls. It can be useful to train volleyball too. In this paper, we presented a comprehensive framework to track players in the indoor volleyball based on combination of Mean Shift, GMM and Kalman Filter. The presented algorithm is new, efficient and rapid so that we can track players in complex situation and frequent occurrence of occlusion. In addition, this method has the capability to recognize and track players when some of them are happening to be at the same line. To test the proposed method, we need a comprehensive dataset of volleyball that consists all the play rules and is not just limited to a specific action. For lacking such a dataset in volleyball, we presented a complete view-dependant volleyball video dataset under the uncontrolled conditions in our previous paper [18]. Therefore, it is possible to test the proposed method on all the images of this dataset, which is prepared by a fixed camera. The accurate obtained results of the players tracking will show us the efficiency of proposed framework. The reminder of this paper is organized as following: The proposed framework to track players in volleyball is presented in section 2. In this section, first we shortly explain the GMM, Mean Shift and Kalman Filter algorithms. Section 3 presents the used dataset for testing the method. The experiment is presented in section 4 and Conclusions are shown in section 5.

II. PROPOSED FRAMEWORK

Nowadays, multi- target tracking is a very active research field, due to its wide practical applicability in video processing. Whenever this problem is considered in sport plays, we encounter more complexity. The most important problem for tracking the players is multi-player occlusion. The general framework presentation to track players in volleyball in a way that solve the occlusion problem well, helps to analyze this play in sports videos. The below section presents the proposed framework for tracking the players in volleyball. As mentioned above, the general framework is a combination of GMM, Mean Shift and Kalman Filter algorithms. At first, we explained the

GMM, Mean Shift and Kalman Filter and then presented the proposed framework.

A. Gaussian Mixture Model

Before we start with tracking of moving objects, we need to extract moving objects from the background. Background subtraction is one of the most common approaches for detecting foreground objects from video sequences. Recently, some statistical methods are used to extract change regions from the background. The Gaussian Mixture Model is the most representative background model [19]. Each background pixel is modeled using a mixture of Gaussian distributions. If the value of a pixel at time t in RGB or some other color space is denoted by $X^{(t)}$, then for Pixel-based background subtraction is needed for distinction of background (BG) and foreground (FG) pixels. Considering Bayesian decision and assuming $P(FG) = P(BG)$, the pixel belongs to the background if:

$$P(X^{(t)}|BG) > C_{thr} \quad (1)$$

Where C_{thr} is a threshold value. Therefore, we use GMM with M components:

$$P(X|BG) = \sum_{m=1}^M \pi_m N(x, \mu_m, \Sigma_m) \quad (2)$$

Where M is the number of distributions, π_m is an estimate of the weight (which portion of the data is calculated for by this Gaussian) of the m^{th} Gaussian in the mixture at time t , μ_m is the mean value of the m^{th} Gaussian in the mixture at time t , Σ_m is the covariance matrix of the m^{th} Gaussian in the mixture at time t , and where N is a Gaussian probability density function. For given a new data sample $x^{(t)}$ at time t , the recursive update equations are presented in [20]. Usually, the intruding foreground objects will be represented by some additional clusters with small weights π_m . Therefore, background model is approximated by the first B largest clusters:

$$P(X|\tilde{x}_T, BG) \sim \sum_{m=1}^B \pi_m N(X; \mu_m, \sigma_m^2 I) \quad (3)$$

If the components are sorted to have descending weights π_m we have:

$$B = \arg \min_b \left(\sum_{m=1}^b \pi_m > (1 - c_f) \right) \quad (4)$$

Where c_f is a measure of the maximum portion of the data that can belong to foreground objects without influencing the background model. All pixels $X^{(t)}$ which do not match any of these components, will be marked as foreground.

B. Mean Shift Algorithm

The Mean Shift algorithm is a nonparametric technique to locate density extremes or modes of a given distribution by an iterative procedure [21, 22]. A target is usually defined by a rectangle or an ellipsoidal region in the image. Most existing target tracking schemes use the color histogram to represent the rectangle or ellipsoidal target. Denote by $\{x_i^*\}_{i=1 \dots n}$ the

normalized pixel positions in the target region, which is supposed to be centered at the origin point. The target model q corresponding to the target region is computed as

$$\hat{\mathbf{q}} = \{ \hat{q}_u \}_{u=1 \dots D} \quad \sum_{u=1}^D \hat{q}_u = 1$$

$$q_u = C \sum_{i=1}^n \mathcal{T}(\|x_i^*\|^2) \Psi[\mathcal{F}(x_i^*) - u] \quad (1)$$

Where q_u represent the probabilities of feature u in target model q , D is the number of feature spaces, Ψ is the Kronecker delta function, $\mathcal{F}(x_i^*)$ associates the pixel x_i^* to the histogram bin, $\mathcal{T}(x)$ is an isotropic kernel profile and constant C is a normalization function defined by

$$C = \frac{1}{\sum_{i=1}^n \mathcal{T}(\|x_i^*\|^2)} \quad (2)$$

Similarly, the target candidate model $\hat{\mathbf{p}}(y)$ corresponding to the candidate region is given by

$$\hat{\mathbf{p}}(y) = \{ \hat{p}_u(y) \}_{u=1 \dots D} \quad \sum_{u=1}^D \hat{p}_u = 1$$

$$p_u(y) = C_h \sum_{i=1}^{n_h} \mathcal{T}(\|\frac{y-x_i}{h}\|^2) \Psi[\mathcal{F}(x_i) - u] \quad (3)$$

$$C_h = \frac{1}{\sum_{i=1}^{n_h} \mathcal{T}(\|\frac{y-x_i}{h}\|^2)} \quad (4)$$

Where $p_u(y)$ represents the probability of feature u in target model $\hat{\mathbf{p}}(y)$, $\{x_i^*\}_{i=1 \dots n_h}$ denote the pixel positions in the target candidate region centered at y , h is the bandwidth and constant C_h , is a normalization function. In order to calculate the likelihood of the target model and the candidate model, a metric based on the Bhattacharyya coefficient is defined between the two normalized histograms $\hat{\mathbf{p}}(y)$ and \mathbf{q} as follows:

$$\hat{\rho}(y) \equiv \rho[\hat{\mathbf{p}}(y) \cdot \mathbf{q}] \quad (5)$$

The distance between $\hat{\mathbf{p}}(y)$ and \mathbf{q} is then defined as

$$d(y) = \sqrt{1 - \rho[\hat{\mathbf{p}}(y) \cdot \mathbf{q}]} \quad (6)$$

Minimizing the distance (6) is equivalent to maximizing the Bhattacharyya coefficient (5). The iterative optimization process is initialized with the target location y_0 in the previous frame. Using Taylor expansion around $p_u(y_0)$, the linear approximation of the Bhattacharyya coefficient (5) is obtained as

$$\rho[\hat{\mathbf{p}}(y) \cdot \mathbf{q}] \approx \frac{1}{2} \sum_{u=1}^D \sqrt{p_u(y_0) q_u} - \frac{C_h}{2} \sum_{i=1}^{n_h} \Delta_i \mathcal{T}(\|\frac{y-x_i}{h}\|^2) \quad (7)$$

Where

$$\Delta_i = \sum_{u=1}^D \sqrt{\frac{q_u}{p_u(y_0)}} \Psi[\mathcal{F}(x_i) - u] \quad (8)$$

Since the first term in (7) is independent of y , to minimize the distance in (6) is to maximize the second term in (7). In the iterative process, the estimated target moves from y to a new position y_1 , which is defined as

$$y_1 = \frac{\sum_{i=1}^{n_h} x_i \Delta_i g(\|\frac{y_0-x_i}{h}\|^2)}{\sum_{i=1}^{n_h} \Delta_i g(\|\frac{y_0-x_i}{h}\|^2)} \quad (9)$$

where $g(x) = -\mathcal{T}'(x)$, assuming that the derivative of $\mathcal{T}(x)$ exists for all $x \in [0, \infty)$, except for a finite set of points.

C. Kalman Filter

A Kalman Filter is used to estimate the state of a linear system where the state is assumed to be distributed by a Gaussian. Kalman Filter consists of five equations and it can be divide them into two groups: time update equations and measurement update equations [23]. The time update equations are responsible for projecting forward (in time) the current state and error covariance estimates to obtain the a priori estimates for the next time step. The measurement update equations are responsible for the feedback—i.e. for incorporating a new measurement into the a priori estimate to obtain an improved a posteriori estimate. The time update equations can also be thought of as predictor equations, while the measurement update equations can be thought of as corrector equations. The state variable and observation of Kalman Filter in this paper are object locations, its velocity and the width of the rectangle which represent the width of a player. The state-space representation of the tracker used in the Kalman Filter is given in (10).

$$\begin{bmatrix} x_{k+1} \\ y_{k+1} \\ x'_{k+1} \\ y'_{k+1} \\ W_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & 0 & \Delta k & 0 & 0 \\ 0 & 1 & 0 & \Delta k & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_k \\ y_k \\ x'_k \\ y'_k \\ W_k \end{bmatrix} + w_k \quad (10)$$

where, x_{k+1} and y_{k+1} are the predicted coordinates of the object and x'_k and y'_k are the velocities in the respective direction, W_k represents the width of the Human rectangle, Δk represents the time interval of state correction and w_k is the white Gaussian noise with diagonal variance Q . The predicted coordinates and dimensions of the rectangle are used to locate the player in the present frame. When the players are distinguished, the Kalman vector is updated using the measurement equation as shown in (11).

$$\begin{bmatrix} \varphi x_{k+1} \\ \varphi y_{k+1} \\ \varphi W_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & 0 & \Delta k & 0 & 0 \\ 0 & 1 & 0 & \Delta k & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_{k-\Delta k} \\ y_{k-\Delta k} \\ x'_{k-\Delta k} \\ y'_{k-\Delta k} \\ W_{k-\Delta k} \end{bmatrix} + v_k \quad (11)$$

Where φx_{k+1} and φy_{k+1} are the measured coordinates, the value φW_{k+1} is the measure width of Human at time $k+1$ and v_k is white Gaussian noise with diagonal variance R . The position, velocity and acceleration are updated based on the values obtained in the present frame and the data from the previous frame.

D. Schematic Description

In this section, we have presented a suitable framework to recognize and tracking the players in volleyball which is illustrated in “Fig. 1”.

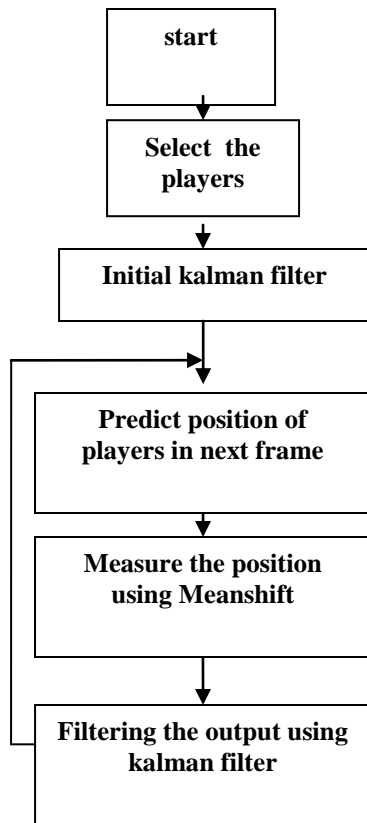


Figure 1. The proposed tracking algorithm scheme.

The presented tracking system for the play is planned in such a way that allow to operator can choose the players who are intended to be tracked. In this way the operator has the possibility to choose 1 to 12 players in the field, so the tracking of specific players will be possible. Using color histogram to track players will be suitable and useful. Mean Shift is rapid and accurate tracking algorithm, which uses color histogram to find the position of the object. Choosing players by operator, there will be found target model q , for every player in Mean Shift. Then window pixels of target candidate being weighted by color histogram of target model and target candidate. After weighing pixels, calculating the center of mass window. Actually, the center of mass the window shows the new place of the object. An important problem in weighing objects is some pixels that are the part of the background gets a high weigh because of the similarity of their color with the object

color and sometimes it lead to fault to find the object. To avoid this problem, we have to enter just the pixels in accounting the center of mass, which are chosen as background by GMM. In such way, more accuracy would be possible to find the new position of player. Calculated the center of mass window until minimize the distance between the new location of player and the predicted position. Finally filtered the newly obtained place and the repetition presented algorithm for every chosen player by operator and for every frame lead to tracking of player in volleyball. Kalman Filter is used to predict the next position of every player in next frame in this work. This position is actually the place of target candidate in Mean Shift algorithm. The place of player calculated by Mean Shift and GMM algorithms then the existing noise will be removed by Kalman Filter.

III. DATASET

To analyze volleyball and assessment of presented method, we need a comprehensive dataset of volleyball that is not specified to a special motion and support all the actions and play rules. The existing dataset such as KHT [24], UCF Sports Action Dataset [25], UCF YouTube Action Dataset [26], IXMAS[27] and Leads Sports Pose Dataset [28] are sets that consists special human motions and sometimes different sports such as volleyball, baseball, soccer, etc. Some of these sets consists images of specific motion in volleyball like spiking and are prepared only to diagnosis the special action and doesn't have a coherency of a full play. In addition, these motions are studied in controlled setting without occlusion. For the frequent occurrence of occlusion in volleyball caused of the excess of players in the proportion with field area, it is necessary to has a dataset of volleyball that consist of complicated occlusion to test the presented method. Lost the comprehensive dataset of volleyball in uncontrolled situation had to us to present a complete view depended volleyball dataset in [18] which called PVD¹. We prepared data in uncontrolled conditions and different viewpoints when planning dataset. The samples of images of dataset sequences are shown in “Fig. 2”.

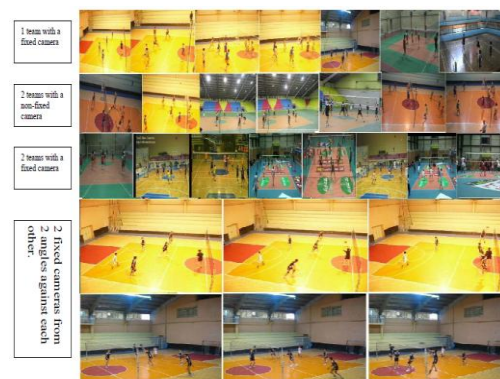


Figure 2. Example of PVD dataset [18].

¹ Persian's Volleyball Dataset

IV. EXPERIMENT

This section explains the obtained results of presented algorithm. We have used the PVD dataset for testing the presented algorithm. This algorithm is applicable for images recorded by a fixed camera. There is possibility analysis a team in volleyball through distinction place of two teams. It is an important because coaches need to study the competitor players for analyze application competition team. One the other hand, the coach needs to study his team players performance in receiving ball, suitable placement and effective passes. There is no need to analyze the competitor team images and complete images of a team and ball with a high quality and maximum zoom works better to analyze. So the prepared images with fixed camera of one team are chosen for testing offer method. Although offer method is applicable for two teams. The results show that the system can correctly recognize players and in some cases the players move in the direction of ball motion and have severe occlusion, the system is able to track the players very well. “Fig. 3” shows results of the presented system with complex situation and severe occlusion.

V. CONCLUSION

The use of computer vision technology in collecting and analyzing statistics during sports matches or training sessions is expected to provide valuable information for tactics improvement. Since in team matches such as volleyball, for the number of players in proportion whit field area, the occlusion happens a lot, it is one of the most challenging in analyzing volleyball. In this paper, we proposed a player tracking method applied for volleyball videos. The proposed framework is based on combination of GMM, Mean Shift and Kalman Filter. In this method, the next position of player is predicted by Kalman Filter and using Mean Shift and GMM algorithms, the exact position of the player is recognized. The testing results on Iran volleyball league official match videos demonstrate that our method can reach high detection and labeling precision, and reliably tracking in cases of scenes such as player occlusion . The method can be utilized for team tactics and player activity analysis in other sports, such as pursuit race, basketball, etc and high-light detection. It also can be applied to surveillance and vision-based human–computer interaction.

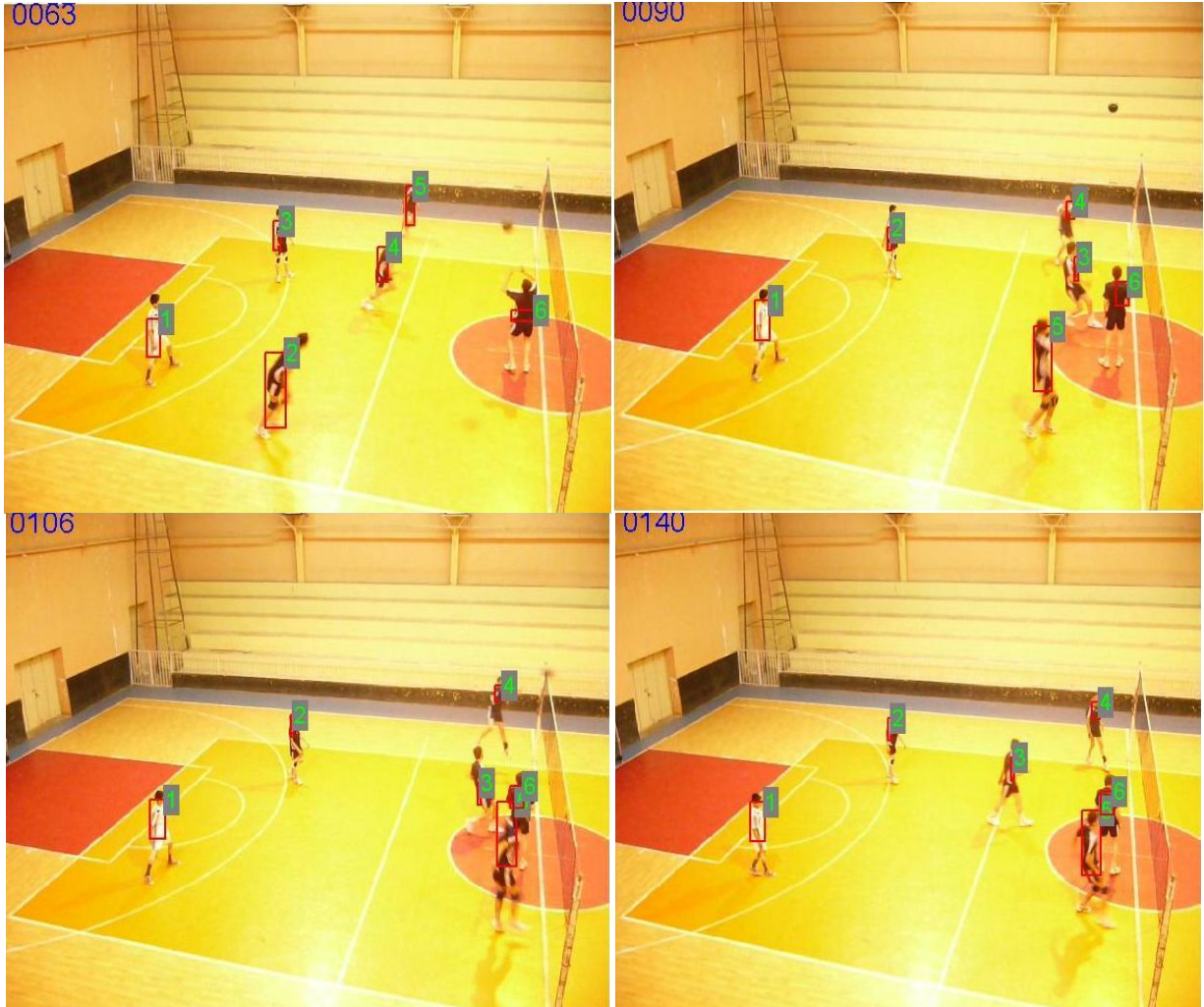


Figure 3. The tracking results of sequences

References

- [1] W S. Erdmann, "Gathering of kinematic data of sport event by televising the whole pitch and track", Proc. 10th Int. Soc. of Biomechanics in Sports Symp. ed R Rodano (Milano, Italy: International Society of Biomechanics in Sports) pp 159–62, 1992.
- [2] C J. Needham, "Doctor's thesis the university of leeds school of computing tracking and modelling of team game interactions", 2003.
- [3] C J. Needham and R D. Boyel, "Tracking multiple sports players through occlusion", congestion and scale 12th British Machine Vision Conf. (Manchester: Manchester University Press) pp 93–102, 2003.
- [4] E H. Khatoonabadi and M. Rahmati, "Automatic soccer players tracking in goal scenes by camera motion elimination", 2003.
- [5] G. Thomas, "Real-time camera tracking using sports pitch markings real-time image process", 2 117–32, 2007.
- [6] J. Liu, X. Tong, W. Li, T. Wang, Y. Zhang and H. Wang, "Automatic player detection, labeling and tracking in broadcast soccer video", 2009.
- [7] S S. Intille, J. Davis and A. Bobick, "Real-time closed-world tracking", Int. Conf. on Computer Vision and Pattern Recognition pp 697–703, 1995.
- [8] J. Pers, M. Bon, S. Kovacic, M. Sibila and B. Dezman, "Observation and analysis of large-scale human motion", Mov. Sci. 21 295–311, 2002.
- [9] M. Kristan, J. Pers, A. Leonardis and S. Kovacic, "A hierarchical dynamic model for tracking in sports, 16th Electrotechnical and computer Science, 2003.
- [10] M. Kristan, M. Pers, S. Kovacic and J. Pers, "Tracking multiple players in sport games using the visual information", Electrotech. Rev. 74 19–24, 2007.
- [11] J. Pers, G. Vuckovic, S. Kovacic and B. Dezman, "A low-cost real-time tracker of live sport events", 2nd Int. Symp. On Image and Signal Processing and Analysis vol 2, pp 362–5, 2001.
- [12] M. Jug, J. Pers, B. Dezman and S. Kovacic, "Trajectory based assessment of coordinated human activity", Proc. 3rd Int. Conf. Computer Vision Systems (Graz, Austria) pp 534–43, 2003.
- [13] Okuma K, "Automatic acquisition of motion trajectories", Master's Thesis vol 2, University of British Columbia, 2003.
- [14] M. Thomas and B. Horst, "A robust multiple object tracking for sport applications", Performance Evaluation for Computer Vision, 2007.
- [15] Th. Mauthner, P.M. Roth and H. Bischof, "Action recognition from a small number of frames", Computer Vision Winter Workshop, 2009.
- [16] H. Salehifar, and A. Bastanfard, "A fast algorithm for detecting, labeling and tracking volleyball players in sports videos", 2011 3rd international conference on signal Acquisition and processing (icasp 2011).
- [17] H. Salehifar, and A. Bastanfard, "Visual tracking of athletes in volleyball sport video", in press.
- [18] H. Salehifar, and A. Bastanfard, "A complete view depended volleyball video dataset under the uncontrolled conditions", in press.
- [19] C. Stauffer and W. Grimson, "Adaptive background mixture models for real-time tracking", In Proceedings CVPR, pp. 246.252, 1999.
- [20] Z. Zivkovic and F. vander Heijden, "Recursive unsupervised learning of finite mixture models". IEEE Trans. on PAMI, vol. 26., no. 5, 2004.
- [21] D. Comaniciu and P. Meer, "Mean shift analysis and applications", In IEEE Int. Conf. Computer Vision, Kerkira, Greece, pp. 1197–1203, 1999.
- [22] C. Beleznaï, B. Frühstück and D. Horst, "Human tracking by fast Mean Shift mode seeking", Journal of Multimedia, vol. 1, no. 1, pp. 1–8, 2006.
- [23] A. Kalman, H. Jazwinski, "Stochastic processes and filtering theory", Academic Press, New York, 1970.
- [24] Ch. Schuld, I. Laptev and B. Caputo, "Recognizing Human Actions: A Local SVM Approach", in Proc. ICPR, Cambridge, UK, 2004.
- [25] M. Sullivan and M. Shah, "Action MACH: A Spatio-temporal maximum average correlation height filter for action recognition", CVPR 2008.
- [26] J. Liu, J. Luo and M. Shah, "Recognizing Realistic Actions from Videos "in the Wild"", IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2009.
- [27] K. Reddy, M. Shah, "Incremental action recognition using feature-tree", Computer Vision, 2009 IEEE 12th International Conference on Computer Vision, 2009.
- [28] S. Johnson and M. Everingham, "Clustered pose and nonlinear appearance models for human pose estimation", In Proceedings of the 21st British Machine Vision Conference, (BMVC2010).