# Data analysis of basketball game performance based on bivariate poisson regression model

## Ke Shen[*]

*Physical Education Institute, Hunan University of Technology, Zhuzhou 412007, Hunan, China*

**Abstract**

Conventional methods used to process two-dimensional discrete data will produce large errors and have a narrow scope of application. Along with the development of mathematical theories and computer technologies, some scholars propose to process two-dimensional discrete data by bivariate Poisson regression model, which takes correlation among data sets into consideration and has excessive variability so that results of data analysis can be more accurate. This paper firstly introduces bivariate Poisson distribution and bivariate Poisson regression model, and then uses this model to analyze performance data of each team in regular seasons of 2013-2014 CBA (China Basketball Association) and 2012-2013 NBA (National Basketball Association), and predict performance in post seasons. Through comparison between actual results and results of double independent Poisson distribution, this model can better predict game performance.

## 1 Introduction

Two-dimensional data is mainly divided into two categories, namely two-dimensional continuous data and two-dimensional discrete data [1]. The former is mainly processed by bivariate continuous distribution. With the development of mathematical theories and computer technologies, people have basically mastered processing methods of two-dimensional continuous data [2]. However, because of calculation complexity, no progress has been made in terms of processing two-dimensional discrete data [3,4]. At present, there are mainly two methods for processing two-dimensional discrete data: one is approximation of data by bivariate continuous distribution. For example, assume that data obeys bivariate normal distribution, and then bivariate normal distribution processing means can be used to analyze data, so that calculation could be simple [5]. The other is to assume that two data sets in two-dimensional discrete data are mutually independent [6]. In this way, two-dimensional discrete data can be converted into two sets of one-dimensional and mutually independent discrete data, and thus the calculation is simplified. For instance, assume that data set $M_1$ and $M_2$ in two-dimensional discrete data sets are mutually independent, and:

$$M_i \sim P\lambda_i, i = 1, 2. \tag{1}$$

After using hypothesis test, two-dimensional discrete data sets can be analyzed according to two independent Poisson distributions, and this method is called double independent Poisson distribution [7].

In basketball games, scores of two teams can be seen as two discrete data sets, and currently basketball scores are mainly analyzed by the above two processing methods. However, these methods have defects: firstly, process scores of discrete basketball games by approximating as continuous distribution, and then it is often impossible to find corresponding continuous distribution, or large errors may be resulted in approximation because of too little amount of data [8]. Secondly, consider that scores of two basketball teams are mutually independent, so as to analyze performance of each team respectively [9]. But in fact, during games, scoring ability, pace and home-away environment of one side will have an influence on the other. Thus, such practice is not rigorous and is very prone to errors [10].

To solve the above defects, Karlis et al proposed to analyze sports game data by bivariate Poisson regression model in 2003. Bivariate Poisson regression model can process two-dimensional discrete data effectively, and bivariate Poisson distribution that it uses was put forward by Holgate in 1964 [11]. During the development, scholars put forward various methods to generate bivariate Poisson distribution probability density function, including bivariate binomial distribution limit acquisition method and trivariate reduction method etc. Compared with conventional processing means, bivariate Poisson regression model takes correlation among data sets into consideration and has excessive variability, so that results of data analysis can be more accurate and have a wider scope of application [12]. Based on the above advantages, bivariate Poisson regression model has been widely used in planning prior insurance rate, prediction of fertility rate and unemployment rate etc. This paper firstly introduces bivariate Poisson distribution and regression model, and then analyzes basketball game performance data by bivariate Poisson regression model.

---

[*] *Corresponding author's* e-mail: shenke111@yeah.net

## 2 Bivariate Poisson distribution and regression model

Trivariate reduction method was proposed by Kocher-lakota in 1992 to calculate probability density function of bivariate Poisson distribution. As this method is widely used, this paper mainly uses trivariate reduction method to analyze and introduce bivariate Poisson distribution and regression model.

Assume that, $X_1$ $X_2$ $X_3$ are random variables that obey Poisson distribution, and are mutually independent with coefficients of $\lambda_1$ $\lambda_2$ and $\lambda_3$ respectively. Build another set of random variables $Y_1$ and $Y_2$, where $Y_1 = X_1 + X_3$ and $Y_2 = X_2 + X_3$, and then joint distribution $(Y_1, Y_2)$ obeys bivariate Poisson distribution. Its probability density function is:

$$P(Y_1 = y_1, Y_2 = y_2) =$$
$$= P(X_1 + X_3 = y_1, X_2 + X_3 = y_2) =$$
$$= \exp\{-(\lambda_1 + \lambda_2 + \lambda_3)\} * \frac{\lambda_1^{y_1}}{y_1!} * \frac{\lambda_2^{y_2}}{y_2!} * \quad . \quad (2)$$
$$* \sum_{i=1}^{\min(y_1, y_2)} \binom{y_1}{i} \binom{y_2}{i} i! (\frac{\lambda_3}{\lambda_1 \lambda_2})^i$$

In addition, covariance is $\text{Cov}(Y_1, Y_2) = \lambda_3$, expected marginal distribution is $EY_1 = \lambda_1 + \lambda_3$ $EY_2 = \lambda_2 + \lambda_3$, and marginal variance is $\text{Var}(Y_1) = \lambda_1 + \lambda_3$, $\text{Var}(Y_2) = \lambda_2 + \lambda_3$. It can be seen that covariance $\text{Cov}(Y_1, Y_2) = \lambda_3$ represents the correlation between two sets of random variables. If $\lambda_3 = 0$ two sets of random variables are mutually independent, and then two-dimensional discrete data set obeys the abovementioned double independent Poisson distribution.

Then, define variable difference as $Z_1 = Y_1 - Y_2$, and probability density function of $Z_1$ is:

$$P(Z_1 = z) = \exp\{-(\lambda_1 + \lambda_2)\} \left(\frac{\lambda_1}{\lambda_2}\right)^{\frac{z}{2}} I_z(2\sqrt{\lambda_2 \lambda_3}). \quad (3)$$

Where, function I is modified Bessel Function. By where, function I is modified Bessel Function. By subtracting probability density function expression of $Z_1$ and probability density functions of two independent Poisson distribution, the same results can be obtained. However, according to the above analysis, each set of random variables in bivariate Poisson distribution is related to $\lambda_3$, which is not 0, and thus probability of variable difference in bivariate Poisson distribution is actually different from variable difference probability of two independent Poisson distributions.

Define the sum of variables as $Z_2 = Y_1 + Y_2$, and then the expected marginal distribution is $EZ_2 = \lambda_1 + \lambda_2 + 2\lambda_3$, and marginal variance is $\text{Var}(Z_2) = \lambda_1 + \lambda_2 + 4\lambda_3$. As $\lambda_3$ is

positive, $\text{Var}(Z_2) > EZ_2$, and this property is called excessive variability.

Building method of bivariate Poisson regression model is similar to that of single-variable Poisson regression model, namely analyze various impact factors of data, introduce to model as covariates and then analyze actual situation. Take $(Y_{1m}, Y_{2m})$ as the $m^{th}$ sample data set, where m=1,2,3……n and each data set obeys bivariate Poisson distribution, introduce parameters $\lambda_{1m}$, $\lambda_{2m}$ and $\lambda_{3m}$ to analysis process as covariates, and bivariate Poisson regression model can be obtained as follows:

$$\begin{cases} (Y_{1m}, Y_{2m}) \sim P(\lambda_{1m}, \lambda_{2m}, \lambda_{3m}) \\ \log(\lambda_{1m}) = a_m \beta_1 \\ \log(\lambda_{2m}) = b_m \beta_2 \\ \log(\lambda_{3m}) = c_m \beta_3 \end{cases} . \quad (4)$$

Where $\beta$ is regression coefficient vector, $a_m$ $b_m$ and $c_m$ are covariate vectors with $m = 1, 2, 3, ......, n$. Values of covariate vectors change according to specific problems, and covariate vectors of parameters $\lambda_{1m}$, $\lambda_{2m}$ can $\lambda_{3m}$ be the same or not. If covariate vector values are the same, it is called consistency hypothesis. As observed values and covariate vectors of samples are known, bivariate Poisson regression model can be obtained by merely re-estimating to obtain value of regression coefficient vector. Here, maximum likelihood estimation is mainly used to estimate regression coefficient, and generally no covariate is introduced to parameter $\lambda_3$ in order to simplify calculation. Then, logarithm of likelihood function can be expressed as:

$$\log L = -\sum_{m=1}^{n}(\lambda_{1m} + \lambda_{2m} + \lambda_3) + \sum_{m=1}^{n}\log \varnothing \quad , \quad (5)$$

$$\varnothing = \sum_{j=0}^{\min(y_{1m}, y_{2m})} \frac{\lambda_{1m}^{y_{1i}-j} \lambda_{2m}^{y_{2i}-j} \lambda_3^{j}}{(y_{1m}-j)!(y_{2m}-j)!j!}. \quad (6)$$

So, the obtained maximum value of $\log L$ is the maximum value of likelihood function, and generally the derivative of function shall be 0 in order to get its maximum value. Therefore, take partial derivatives of parameters $\beta_1$, $\beta_2$ and $\lambda_3$ b $\log L$ respectively, and:

$$\begin{cases} \frac{\partial \log L}{\partial \beta_1} = \sum_{m=1}^{n} \lambda_{1m} u_m (P_{10} - 1) = 0 \\ \frac{\partial \log L}{\partial \beta_2} = \sum_{m=1}^{n} \lambda_{2m} v_m (P_{01} - 1) = 0 \\ \frac{\partial \log L}{\partial \lambda_3} = -n + \sum_{m=1}^{n} P_{11} = 0 \\ P_{jk} = \frac{\Phi(y_{1m} - j, y_{2m} - k)}{\Phi(y_{1m}, y_{2m})} \end{cases} . \quad (7)$$

## 3 Application of bivariate Poisson regression model in basketball game data analysis

Cummins et al proposed in 1983 that bivariate regression model can be applied to insurance claims and clarification of insurance premium etc, which was the earliest application of this model. Compared with conventional double independent Poisson distribution, results obtained through this model are more consistent with actual situation, and its application also obtains good benefits. With the promotion of this model, Karlis et al proposed in 2003 to analyze sports game data and predict results by bivariate Poisson regression model. Karlis also compared prediction results of football and water polo games by bivariate Poisson regression model, prediction results of double independent Poisson distribution and actual situation, and results showed that results of bivariate Poisson regression model were more accurate.

Generally speaking, performance data of two teams in sports games is conventionally analyzed by double independent Poisson distribution. In this way, it is believed that score distribution of two teams is mutually independent and obeys Poisson distribution. However, in fact, score distribution is correlated and scoring ability, pace and home-away environment of one side will have an influence on the other. In particular, during basketball game with fast pace, pace of one side will inevitably affect the other if scores rise alternately. Therefore, game performance data shall be analyzed and predicted by bivariate Poisson regression model, which has two advantages: firstly, correlation factors in actual situation are included into the model, and secondly, double independent Poisson distribution does not consider about this situation as actually collected data is of excessive variability so that expectation is the same with variance. Bivariate Poisson regression model is featured by excessive variability and thus can better meet requirements of data processing. Thus, the author uses bivariate Poisson regression model to analyze data of regular seasons of 2013-2014 CBA (China Basketball Association) and 2012-2013 NBA (National Basketball Association) respectively, and predict performance in post seasons.

### 3.1 CBA GAME PERFORMANCE DATA ANALYSIS BY BIVARIATE POISSON REGRESSION MODEL

At present, general point rules of basketball games are 2 scores for winning one game, 1 score for losing one game and 0 score for giving up. Table 1 shows points of each team in 34 rounds of CBA regular season, and specific score data can be seen on official website of China Basketball Association.

TABLE 1  Points of each team in 34 rounds of 2013-2014 CBA regular seasons

| Ranking | Team | Winning | Losing | Points | Average scores | Average lost scores |
|---|---|---|---|---|---|---|
| 1 | Guangdong Winnerway | 30 | 4 | 64 | 100.3 | 88.6 |
| 2 | Xinjiang Guanghui | 26 | 8 | 60 | 104.3 | 93.5 |
| 3 | Dongguan Men's Basketball | 25 | 9 | 59 | 105.5 | 100.7 |
| 4 | Beijing Jinyu | 23 | 11 | 57 | 105 | 98.6 |
| 5 | Zhejiang Guangsha | 21 | 13 | 54 | 108.6 | 105.7 |
| 6 | Tianjin Reapal | 20 | 14 | 54 | 104.3 | 103.6 |
| 7 | Liaoning Hengye | 20 | 14 | 54 | 101 | 99.6 |
| 8 | Shanghai Men's Basketball | 20 | 14 | 54 | 98.5 | 96.2 |
| 9 | Shandong Gold | 19 | 15 | 53 | 94.6 | 92.5 |
| 10 | Fujian Men's Basketball | 16 | 18 | 50 | 107.8 | 108.4 |
| 11 | Jiangsu China Railway | 15 | 19 | 49 | 99.2 | 100.3 |
| 12 | Sichuan Aijia | 14 | 20 | 48 | 97.5 | 105.7 |
| 13 | Zhejiang Chouzhou | 13 | 21 | 47 | 106.8 | 106.8 |
| 14 | Jilin Rural Commercial Bank | 12 | 22 | 46 | 101 | 105.9 |
| 15 | Foshan Rural Commercial Bank | 11 | 23 | 45 | 100.5 | 105.8 |
| 16 | Shanxi Fenjiu | 10 | 24 | 44 | 101.5 | 104.2 |
| 17 | Bayi Shuanglu | 6 | 28 | 40 | 92.7 | 101.3 |
| 18 | Qingdao Double Star | 5 | 29 | 39 | 102.8 | 114.6 |

Assume that points of one side X is $X_i$, and the other is $Y_i$, it can be seen from formula (7) that points of each game meet the following distribution:

$$\begin{cases} (X_i, Y_i) \sim P(\lambda_{1i}, \lambda_{2i}, \lambda_{3i}) \ i = 1, 2, \ldots n \\ \log(\lambda_{1i}) = \mu + h + atth + defg \\ \log(\lambda_{2i}) = \mu + attg + defh \end{cases} , \qquad (8)$$

where, n is the total number of rounds, $\mu$ is a constant, h denotes the influence of home field factors on the game, atth and defh denote attach and defense scores of home team, while attg and defg denote attach and defense scores of away team. In order to make model parameters more identifiable, it is necessary to use standard constraint con-

ditions. Thus, when selecting values of $\lambda_1$ and $\lambda_2$, consider that two teams play on a neutral field and have the same physical strength. In addition, parameters of attach and defense shall be those that represent average attach and defense abilities of teams. Parameter $\lambda_3$ denotes random factors that affect game results, such as pace and audience atmosphere etc. Thus, covariance $\lambda_{3i}$ can be expressed as:

$$\log(\lambda_{3i}) = \beta^{con} + \gamma_1 \beta_{h_i}^{home} + \gamma_2 \beta_{g_i}^{away} , \qquad (9)$$

where, $\beta^{con}$ is a constant, $\beta_{hi}^{home}$ and $\beta_{gi}^{away}$ denote parameters of home and away field abilities respectively, $\gamma_1$ and $\gamma_2$ are bivariate indicating parameters with values of 0 or 1

according to actual conditions. If $\gamma_1 = \gamma_2 = 0$, covariance is invariable, and if $\gamma_1 = 1$ and $\gamma_2 = 0$, covariance is only related to home team. Expected scores of two teams are:

$$\begin{cases} E(X_i) = (1-p)(\lambda_{1i} + \lambda_{2i}) + p\theta_1 \\ E(Y_i) = (1-p)(\lambda_{2i} + \lambda_{3i}) + p\theta_1 \end{cases}, \qquad (10)$$

where, $P$ and $\theta_1$ are estimated mixing ratio and expansion factor respectively. According to statistics of existing game data, the following game performance can be estimated by formula (10). Difficulty of solution-finding by this model lies in estimation of regression coefficient vector $P$, and commonly used estimation methods include Newton-Raphson Method and Expectation Maximization Algorithm. Thus, this paper will analyze game performance data by these two methods respectively.

Firstly, solve bivariate Poisson regression model of CBA game performance data by Newton-Raphson Algorithm, iterative formula of which is:

$$x^{k+1} = x^k - F'(x^k)^{-1} F(x^k), \text{ k=0,1,2...,}$$

where $F(x^k)$ is a functional matrix. When solving formula (9) by this iterative formula, take second-order derivative of equation in formula (9). Similarly, from the initial value, estimated value of regression coefficient in the model can be obtained through repeated iteration.

Newton-Raphson Algorithm has fast rate of convergence, but has a high requirement for selection of initial value and requires that the selected initial value shall be as close as possible to accurate value of function, so that results can converge, otherwise, calculation effects cannot be reached. When determining initial value, value estimated under independent conditions is generally selected.

By solving the model, the author analyzes the game results, predicts results of post season, compares with actual situation and shows the results in table 2. It shall be noted that winning-losing relationship between teams is converted into points for the convenience of calculation. For example, team A wins 3 out of 5 rounds with team B, and the point is 8-7.

TABLE 2    Actual and predicted results of 2013-2014 CBA post season (take $\mu$=0)

| Round | Both sides | Actual winning-losing | Actual points | Predicted points |
|---|---|---|---|---|
| 1/4 Final | Guangdong Winnerway-Shanghai Men's Basketball | 3-0 | 6-3 | 5.2-3.4 |
| 1/4 Final | Xinjiang Guanghui-Liaoning Hengye | 3-1 | 7-5 | 6.3-3.9 |
| 1/4 Final | Dongguan Men's Basketball-Tianjin Reapal | 3-1 | 7-5 | 6.1-4.1 |
| 1/4 Final | Beijing Jinyu-Zhejiang Guangsha | 3-1 | 7-5 | 5.9-4.3 |
| Semifinal | Guangdong Winnerway-Beijing Jinyu | 2-3 | 7-8 | 7.5-6 |
| Semifinal | Xinjiang Guanghui-Dongguan Men's Basketball | 3-0 | 6-3 | 4.5-3 |
| Final | Xinjiang Guanghui-Beijing Jinyu | 2-4 | 8-10 | 8.4-7.2 |

From comparison with actual points, predicted points can predict winning-losing relationship in games well, except deviation of prediction of Beijing Jinyu.

### 3.2 NBA GAME PERFORMANCE DATA ANALYSIS BY BIVARIATE POISSON REGRESSION MODEL

NBA has 30 teams, each of which shall play 82 games in regular seasons, and thus more data can be accumulated compared with games in CBA. Here, the author uses expectation maximization algorithm to solve bivariate Poisson regression model of NBA game performance data.

Expectation maximization algorithm contains generally two steps: firstly, calculate expectation, namely estimate expected values of unknown parameters and give current parameter estimation. Secondly, maximize expectation, namely re-estimate distribution parameters to achieve maximum data likelihood and give estimated expectation of location variable. Main idea of this algorithm is: assume that $f(\theta|E)$ is a posterior density function of data E and $\theta$, $f(\theta|E,F)$ is a posterior density function of data E and F and $\theta$. These two functions are called observation posterior density and addition posterior density respectively. Besides, take $f(G|\theta,E)$ as a conditional density function under

the situation that $\theta$ and data E are determined. Mean value of observation posterior density shall be calculated by expectation maximization method. Then, steps of expectation calculation and maximization are:

$$Q(\theta|\theta^k, E) = \int \log[f(\theta|E,F)]f(G|\theta^k, E)dF, \qquad (11)$$

$$Q(\theta^{k+1}|\theta^k, E) = \max Q(\theta|\theta^k, E), \qquad (12)$$

where, $\theta^k$ and $\theta^{k+1}$ are approximate values after k and k+1 times of iteration respectively. Use formula (11) and (12) for repeated iteration until $|\theta^{k+1} - \theta^k|$ is small enough. For bivariate Poisson regression model, covariates of parameters are also different, and only iterative method can be used to calculate. In expectation maximization method, firstly, calculate conditional expectation of parameters to be estimated; and then maximize the obtained expectation to complete estimation of regression coefficient.

The author analyzes data of 2012-2013 NBA regular season, and specific score data can be seen on official website of NBA and thus will not be given in detail here. Besides, the author also predicts results of post season, compares actual situation and shows results in table 3.

TABLE 3    Actual and predicted results of 2012-2013 NBA post season (take $\mu$ =0)

| Round | Both sides | Actual winning-losing | Actual points | Predicted points |
|---|---|---|---|---|
| Western 1/4 Final | Thunder-Rockets | 4-2 | 10-8 | 8.4-5.8 |
| Western 1/4 Final | Spurs-Lakers | 4-0 | 8-4 | 7.6-4.2 |
| Western 1/4 Final | Nuggets-Warriors | 2-4 | 8-10 | 8.2-6.9 |
| Western 1/4 Final | Clippers-Grizzlies | 2-4 | 8-10 | 7.8-8.3 |
| Eastern 1/4 Final | Heat-Bucks | 4-0 | 8-4 | 7.7-4.1 |
| Eastern 1/4 Final | Knicks-Celtics | 4-2 | 10-8 | 9.7-7.2 |
| Eastern 1/4 Final | Pacers-Hawks | 4-2 | 10-8 | 10.2-7.3 |
| Eastern 1/4 Final | Nets-Bulls | 3-4 | 10-11 | 10.6-10.2 |
| Western Semifinal | Thunder-Grizzlies | 1-4 | 6-9 | 9.3-7.6 |
| Western Semifinal | Spurs-Warriors | 4-2 | 10-8 | 10.3-8.6 |
| Eastern Semifinal | Heat-Bulls | 4-1 | 9-6 | 10.3-7.2 |
| Eastern Semifinal | Pacers-Knicks | 4-2 | 10-8 | 9.7-6.3 |
| Western Final | Spurs-Grizzlies | 4-0 | 8-4 | 9.5-8.4 |
| Eastern Final | Heat-Pacers | 4-3 | 11-10 | 10.6-9.5 |
| Final | Heat-Spurs | 4-3 | 11-10 | 9.6-10.4 |

It can be seen from comparison with actual points that predicted results are basically consistent with actual results. From analysis of the above two cases, results of bivariate Poisson regression model by both Newton-Raphson Method and Expectation Maximization Algorithm can be used to analyze basketball game performance data effecttively and predict results.

### 3.3 COMPARISON BETWEEN BIVARIATE POISSON REGRESSION MODEL AND DOUBLE INDEPENDENT POISSON DISTRIBUTION

For analysis of basketball game data, the greatest differrence between double independent Poisson distribution and bivariate Poisson regression model lies in lack of consideration of score correlation of both teams, namely covariance $\lambda_{3i}$ is equal to 0 in the model.
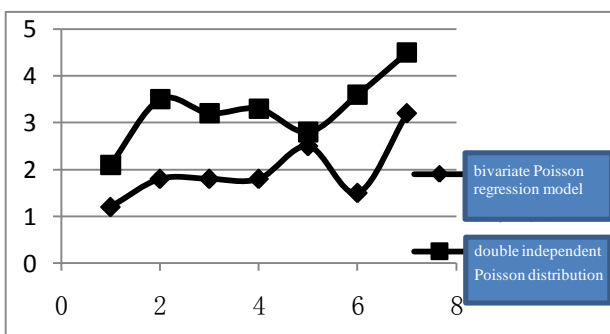


FIGURE 1 Difference between predicted and actual results of bivariate Poisson regression model and double independent Poisson distribution

To verify the superiority of bivariate Poisson regression model, the author analyzes performance of CBA regular season shown in 3.1 and predicts results of post season by double independent Poisson distribution.

The author compares the obtained results with those obtained by bivariate Poisson regression model, and predicts results of a total of 7 post seasons, measures accuracy of prediction by absolute value of difference between predicted points and actual points, and shows results in figure 1. It can be seen that bivariate Poisson regression model has more accurate prediction than double independent Poisson distribution.

### 4 Conclusions

Two-dimensional discrete data can be analyzed by bivariate Poisson regression model, which considers about correlation among data sets and has excessive variability, and thus results of data analysis can be more real and accurate. Based on analysis of bivariate Poisson distribution and bivariate Poisson regression model, this paper builds a model for basketball game performance data, and uses this model to analyze data of 2013-2014 CBA regular season and 2012-2013 NBA regular season and predict results of post seasons. Compared with double independent Poisson distribution, this model has predicted results that are more consistent with actual results, indicating that it is feasible to analyze basketball game performance data by bivariate Poisson regression model.

### References

[1] Fromm E G, RG Quinn. (1989) An experiment to enhance the educational experience of engineering students. *Engineering Education*, **79**(3), 424-429.
[2] Yeomans S R & A Atrens. (2001) A methodology for discipline-specific curriculum development. *International Journal of Engineering Education*, **17**(6), 518-524.

[3] Ashish Kumar Parashar, Rinku Parashar. (2012) Innovations and Curriculum Development for Engineering Education and Research in India. *Procedia-Social and Behavioral Sciences*, **56**(1), 685 – 690.
[4] Behrooz Parhami. (1986) Computer science and engineering education in a developing country: the case of Iran. *Education and Computing*, **2**(4), 231-242.

[5]  Alan B Forsythe, James R Freed, Harvey S Frey. (1975) Programmed Instruction Nucleus (PIN): A Simplified Author-language for Computer-aided Instruction. *Computers in Biology and Medicine*, **5**(6), 77-88

[6]  Hiroki Okubo, Mont Hubbard. (2010) Identification of basketball parameters for a simulation model. *Procedia Engineering*, **2**(2), 3281-3286

[7]  Fazhi Sun. (2013) Comprehensive Evaluation and Research on Teaching Abilities of Basketball to the Normal Colleges Students Specializing in Basketball of Physical Education. *Journal of Digital Content Technology and its Applications*, **7**(2), 379-385.

[8]  Ma Yong. (2013) Study on Evaluation of Basketball Training Effect Based on Fuzzy Comprehensive Evaluation Method. *Advances in Information Sciences and Service Sciences*, **5**(5), 251-258.

[9]  Jiayi Zhu. (2013) Triangular Fuzzy Synthetic Evaluation of Basketball Guard's Offensive Ability. *Advances in Information Sciences and Service Sciences,* **5**(5), 859-865.

[10] CHEN Yongxin. (2013) Comparative Analysis of the Plyometric Training and Dynamic Tensile on the Influence of Male Basketball Players Vertical Jump and Agility. *Journal of Convergence Information Technology*, **8**(10), 1084-1091.

[11] Dong Lifeng. (2013) Utility Analysis of the Effect of Continuous Training on Basketball Players' Physical Improvement. *Journal of Convergence Information Technology*, **8**(6), 663-670.

[12] Hongqiang Li. (2013) The Application of Modern Information Technology in Junior Basketball Training with Linguistic Information. *Journal of Digital Content Technology and its Applications*, **7**(4), 850 -855.

## Authors

**Ke Shen, 16. 4. 1983, China**

**Current position, grades:** A teacher at Physical Education Institute, Hunan University of Technology, China.
**University studies:** Bachelor Degree of Science in Education in 2005 and Master degree of Science in Education from Central China Normal University, China in 2008.
**Scientific interest:** Physical education and physical exercise.