# PREDICTION OF EDUCATIONAL DATA MINING BY MEANS OF A POSTPROCESSOR TOOL

## Oktay Kir, Irina Zheliazkova

**Abstract**:  *A methodology for application of several linear methods of prediction of correct, missing, and wrong knowledge using a postprocessor tool is presented.*

**Keywords**: *teacher, learner, moving average methods, error of prediction, methodology, prediction skill, sessionscript*

 **ACM Classification Keywords**: *Computer and Information Science Education, Knowledge Representation*

## Introduction

In a recent exhaustive survey of *Romero Cr., Ventura S., 2010* [3] prediction has been pointed out as one of the oldest *Educational Data Mining* (*EDM*) task. The variables more often predicted are the learner's performance, knowledge, scores, and mark and the techniques most commonly applied are neural and Bayesian networks, rule-based systems, regression, and correlation analysis.

In a previous authors' paper [2] a teacher's tool for the *EDM* called postprocessor was reported from design, implementation, and user's points of view. The term "postprocessor" stands for processing standardized output data sets after each session, e.g. test, lecture or exercise from a corresponding task-oriented environment. For ensuring the tool's intelligence and its adaptation to the teacher a power and expressive script language called *SessionScript* was implemented. Programming of descriptive statistics, visualization, and correlation analysis techniques was demonstrated using two output data sets respectively from environment for knowledge testing and for exercise task performing.

In business modeling the classical linear and non-linear methods of prediction are referred to as a temporal data mining technique for estimation of unknown values of an observed variable [5]. According to [1] the base line, e.g. time seria with the observed numerical, continuous or discrete values has to meet four important requirements:
a) The results of observations have to be sorted from the earliest to the last one; b) All time periods have to be of equal length; c) The observations have to be fixed at the same time in each period; d) Missing even a single observation is not desirable and missing data has to be complete with estimated ones. If a given base line does not meet any of these requirements it is likely the prediction error to be unacceptable.

The computationally complex techniques listed in the above-mentioned survey are proved as successful for the tasks concerning mainly mediate- and long-term prediction. In this paper studying  simple linear methods for the *L's* short-term prediction of correct, missing, and wrong knowledge of testing is reported using the postprocessor. Firstly the pedagogical experiment carried out for gathering the input data set is described.  The focus of the

paper is on implementation of four methods called moving average in three cases, e.g. correct, missing, and wrong knowledge of testing. A comparative analysis of the corresponding errors of prediction by means of the tool is made to choose the most precise method. Conclusions summarize the methodology proposed for the application of the considered methods using the postprocessor.

## Pedagogical Experiment Description

The data set for implementation of linear methods in the postprocessor for short-term prediction was gathered in the framework of an experiment carried out the academic 2008/2009 year. Four groups of bachelor degree regular students (3-rd year, 1–st semester) specialty "Computer Systems and Technologies" at Rousse University were involved in it (63 students in total). The test session was carried out within the framework of the course in Software Engineering and covered 30 hours taught lecture material. The test was created by the lecturer of the course as an intelligent posttest in order to evaluate the correct, missing, and wrong knowledge, as well as the time undertaken for the test performance.

Each student was accessible through a common device to a template Microsoft Word document. The number of questions was 30 with total scores $P_{max}$ = 352 and planned time $T_{max}$ = 120 min. The questions types were four, namely: multiple choice, unordered keywords, ordered keywords, and unordered pairs [4] and answering was reducing to filling an empty edit field in correspondence with a simple syntax. Depending on its type each question brought different number of scores $p_{\max j}$. A question "*no*" answer or subanswer was interpreted as missing knowledge, and incorrect answers or subanswers as wrong knowledge. After the test performance the student had to upload the fulfilled document back on the common device. The students were also told that the time for the test performance actually is unlimited and together with wrong and missing knowledge will be used as assessment indicators only for research purpose. Later the lecturer manually calculated the questions correct, missing, and wrong knowledge, their total scores  for each student and his/her final mark in the traditional six-based scale: $0 \leq P \leq 0.4* P_{max}$ – "2"; $0.4* P_{max} < P \leq 0.55* P_{max}$ – "3"; $0.55* P_{max} < P \leq 0.70* P_{max}$ – "4"; $0.70* P_{max} < P \leq 0.85* P_{max}$ – "5"; $0.85* P_{max} < P \leq 1.0* P_{max}$ – "6". The experience accumulated during the last decade by Zheliazkova's research group has pointed out that such non-linear scale is acceptable by both teachers and students [6]. The Word document of a "good" student, e.g. received test mark "4" is shown on fig. 1. This experiment confirmed the fact that the Ls go to such intelligent test performance only if they were preliminary self assessed at least with the mark "3". The students also were well motivated and stated that were waiting for an objective and precise test assessment. As a result a tendency of shifting the average test mark from "good" to "very good" also was monitored.

## Tool's Description

In order to solve a new task a new user has to familiar with the data mining techniques, tasks classification, as well as with SESSIONSCRIPT language. The full specification of its first version can be found in [2]. By means of a standard text editor the user can review the script of the programs for related tasks.

A new free-text formulated problem has to be clear, precise, and compact. Although the problem solving is presented as a sequence of steps in practice some steps can be omitted others repeated or interpreted as subproblem solving. To perform each step from the technological scheme the teacher has to know the syntax

and semantics of the corresponding group of commands. In order to enhance the SESSIONSCRIPT language learning the following color coding scheme has been accepted: the correct commands names and symbols for operations in *Aqua*; table, row and column names in *Yellow*; values in *Pink*; unknown keywords and current values of program variables in *Dark grey*; and messages in *Grey*.

РУ "Ангел Кънчев"
катедра "Компютърни Системи и Технологии"
**Многоцелеви групов тест върху лекционния материал**
**по дисциплината "СОФТУЕРНО ИНЖЕНЕРСТВО" (СИ) за студенти редовно обучение 2008/2009**
*Авторски колектив:*
1. Ангел Иванов, студент-бакалавър, 3-ти курс, спец. КСТ
2. Юсуф Хасанов, студент-бакалавър, 3-ти курс, спец. КСТ
3. Юмер Юмер, студент-бакалавър, 3-ти курс, спец. КСТ
4. Денис Ислям, студент-бакалавър, 3-ти курс,спец. КСТ
5. доц. д-р Ирина Желязкова, преподавател по СИ, кат. КСТ
*Цели:*
1. Оценка на изходното ниво знанията;
2. Диагностика на пропуските в лекционния материал;
3. Оценка качеството на теста.
*Скала за оценка:*
от 000 до 127 т. – 2
от 128 до 174 т. – 3
от 175 до 221 т. – 4
от 222 до 268 т. – 5
от 269 до 315 т. – 6
*Очаквано време за изпълнение*: 120 мин
Студент: . Ана Данчева Панова **Фак. номер**:063178 **Група** 20 б

**Тема 1: ПРОЦЕДУРЕН ПОДХОД**　　　　**7 въпроса**　　　　　　**65 т.**

**ВЪПРОС 1:** *Подредени ключови думи*
**Попълнете липсващите ключови думи:** "*Различават два вида ... между модулите: ... и вътрешна. Една програма трябва да се разбие на ... така, че да съществува ... връзка между отделните модули и обратно -... връзки в модулите.*"
**Отговор:** връзки > външни >модули>слаби > силни
**Параметри:** L=2;$Q_t$=10;$C_p$=0.50
**Препратка:** *1.1. Същност на процедурния подход*
**Знания:** 8,0,2
...

**ВЪПРОС 18:** *Неподредени двойки*
**Укажете съответствието между: 1) INTERFACE, 2) IMPLEMENTATION** *на статична библиотека с име StatDLL за нуждите на програмата от фиг. 5 и други програми, използващи ПП от една и съща DLL .*
**Отговор:** 1>a; 1>c ;2>b ;2>d
**Параметри:** L=2;$Q_t$=12;$C_p$=0.66
**Препратка:** *3.4. Създаване на интерфейсен модул за DLL*
**Знания:** 10,0,0
...

**ВЪПРОС 31:** Считате ли, че получената оценка е обективна и точна? Да/Не
**ВЪПРОС 32:** За колко минути попълнихте теста?
**Резултат:** 218,66,31,106,4

**РЕЧНИК**

| | | | | |
|---|---|---|---|---|
| библиотеки | инициализиращите | портове | AStud | MEML |
| библиотеката | интерфейсната | празна | Begin | new(PPrezident) |
| библиотеките | клавиатура | прекомпилирани | Build | PORTW |
| връзките | код | принтер | call | program |
| външни | константи | програмата | const AStud: TStud = | self |
| вътрешна | масиви | програмните | End. | Total |
| главната | методите | променливи | export | Total:=0; |
| глобални | минимална | реализационната | exports | TPrezident |
| датата | модул | свързаност | External | ^TPrezident; |
| двубайтови | модули | сегмент | FAR | Unit  StatDLL; |
| деклариране | нови | силни | index | USES |
| декларацията | обекта | статична | inline | Value : word; |
| диаграми | обектен тип | стека | library | var |
| директиви | обектна | съществуващи | Make | VAR Total : word; |
| дисков файл | отместване | текущия | MEM | with |
| еднобайтови | периферните | фрагменти | MEMW | {$F+} |
| запис | полетата | APrezident^ | | |

Дата: 28.01.09　　　　　　　　ЖЕЛАЕМ ВИ УСПЕХ!

Fig. 1. The word document of a student's test

The left side window (fig. 2) contains the script code of the program written by the teacher, and the right side window the table with the input data set. The description of the problem solved in a free text format can be seen switching from "*Table*" to the "*Description*" tabs (not shown here). The teacher is also recommended to save it in a standard text file serving as a common catalogue of the problems already solved by the members of the course team.

To facilitate the visualization the transposition of the input table is recommended. Each row of the table on fig. 2 corresponds to a test question from 1 to 30, and its first three columns the input data set, e.g. respectively the base line with correct (RA), missing (EA), and wrong (WA) knowledge. The next three columns contain their normalization values (P_RA), (P_EA), and (P_WA) respectively to the maximal scores $p_{\max j}$ the corresponding question. The teacher can choose also the menu-command Window|Variables to view the names and values of the system variables and the program variables. The visualization allows choosing different kinds of diagrams, such as bar, pie, line, and point viewed in separated windows. In order to view the corresponding diagram the table name has to be chosen from the menu-item View|Bars. For the needs of this kind of tasks a power command for visualization of a family of lines in a common coordinate system had been implemented in the postprocessor: PLOT($x_1,x_2,\ldots x_n$, 'propertyName$_1$ = propertyValue$_1$'; $y_1,y_2,\ldots y_n$, 'propertyName$_2$ = propertyValue'; …) propertyName and propertyValue => { TITLE = "<string>", LINESTYLE = DOT | SOLID, LINEWIDTH = <integer>, MARKVISIBLE = TRUE | FALSE, POINTVISIBLE = TRUE | FALSE, POINTSIZE = <integer>, POINTSTYLE = SQUARE | TRIANGLE | CIRCLE | DOWNTRIANGLE | CROSS | DIAGONAL | STAR | DIAMOND, COLOR = CL<COLOR_NAME>}



Fig. 2. The tool's screen with the script and table windows

## Correct, Wrong, and Missing Knowledge Prediction

Hereinafter three linear methods called Moving Average (MA) are remained from the reviewed INTERNET literature [7,8,9,10] where the term *MA* stands for the mean value for a certain period of time.

➢  *Simple Moving Average* (*SMA*) uses average demand for a fixed sequence of periods and is good for a stable demand with no pronounced behavioral samples. The calculated formula for the step averaging procedure is $F_{t+1} = \sum\limits_{i=t-N+1}^{t} F_t$ , where $F_{t+1}$ is the prediction for the $(t+1)^{th}$ period of time; $F_t$ is the actual value at the $t^{th}$ period; *N* is the number of observed periods. The main disadvantage of the method is loosing several predicted values which number is equal to the number of the time periods for the MA.

➢ *Weighted Moving Average* (*WMA*) allows placing a greater emphasis on more recent data in order to reflect changes in demand samples. The weights used are based on the experience of the human predictor. In practice the weighting factors are often chosen to give more weight to the most recent data in the time series and less weight to older one. The corresponding formula is $F_{t+1} = \sum\limits_{i=t-N+1}^{t} w_i F_i$ , $\{w_1, w_2,...,w_n\}$ the vector of weights such that $\sum\limits_{i=1}^{N} w_i = 1$ . Note, that this method does not avoid the disadvantage of the *SMA* and requires a more complex calculation at each step of averaging procedure. Additionally, if the data from each step are not available for analysis, it can be difficult if impossible to reconstruct a changing signal accurately. However, if the number of missing steps is known, the weighting of values in the *MA* can be adjusted to give equal weight to all missing samples to avoid this issue.

➢ *Exponential Moving Average* (*EMA*) is a better trend indicator than the *SMA* as it puts greater weight to most recent data than older ones. Unlike *SMA* and *WMA* the older data never goes away in the calculation of *EMA*. In this case the step calculation formula has the following form $F_{t+1} = (1-\alpha)F_t + \alpha Y_t = F_t + \alpha(Y_t - F_t), t > 1$; $\alpha = 2/(M+1)$; $0 < \alpha < 1$, where the coefficient $\alpha$ is called *smoothing factor* and *M* is called *length*.

The results of applying the above-mentioned methods of prediction on the first base line, e.g. for the correct knowledge are shown on fig. 3. Hereinafter the following color scheme is used for visualization: the observation values of the base line in red, line with prediction values for the *SMA* method (5 periods of time) in green, for the same method but with 15 periods in blue, line with prediction values for the *WMA* method with 5 periods in black, and line with prediction values for the *EMA* method with 5 periods in pink. For the three cases, e.g. correct, wrong, and missing answers the weights for the WMA method, e.g. were calculated as follows:  W1 = 1 / (5+7*1.61803398), W2 = W1*61803398, W3 = W1 + W2, W4 = W2 + W3, W5 = W3 + W4, where Wi is the weight of the ith period (i = 1,…,5). The *WMA* attaches more value to the latest data.
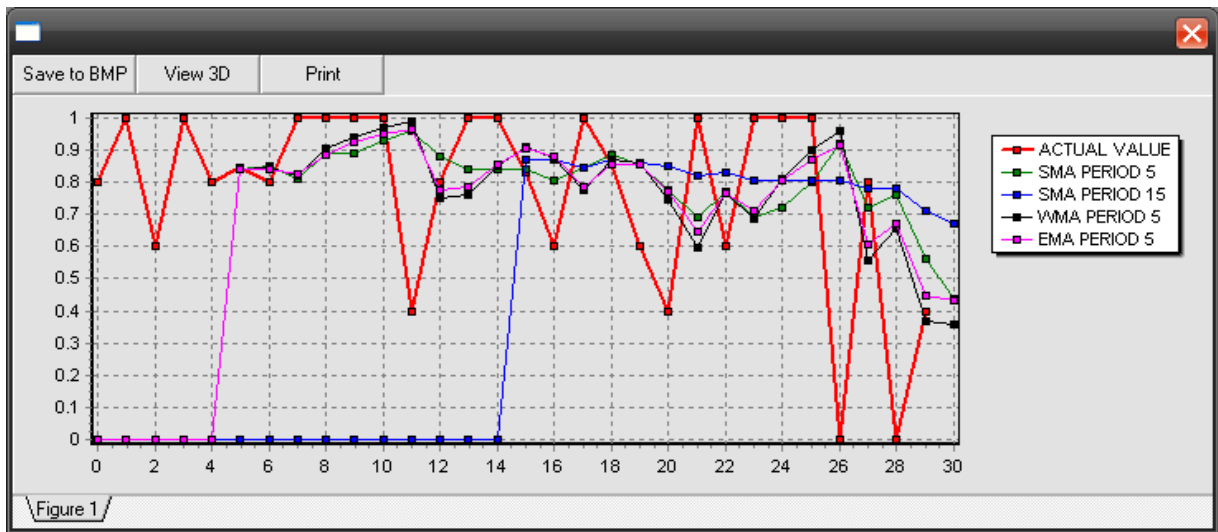
Fig. 3. The window for visualization of three prediction methods together with the base line for correct knowledge

On fig. 4 the results of programming the same four methods of prediction for the wrong knowledge are shown.
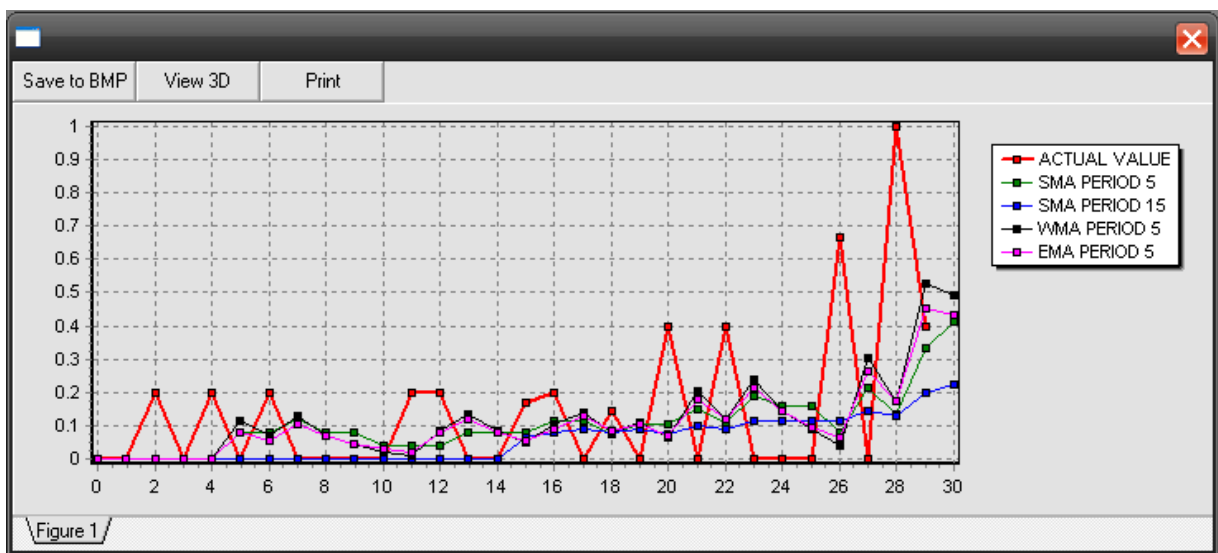


Fig. 4. The window for visualization of three prediction methods together with the base line for wrong knowledge

The results of programming the above-mentioned methods of prediction from the base line for the missing knowledge are shown on fig. 5.
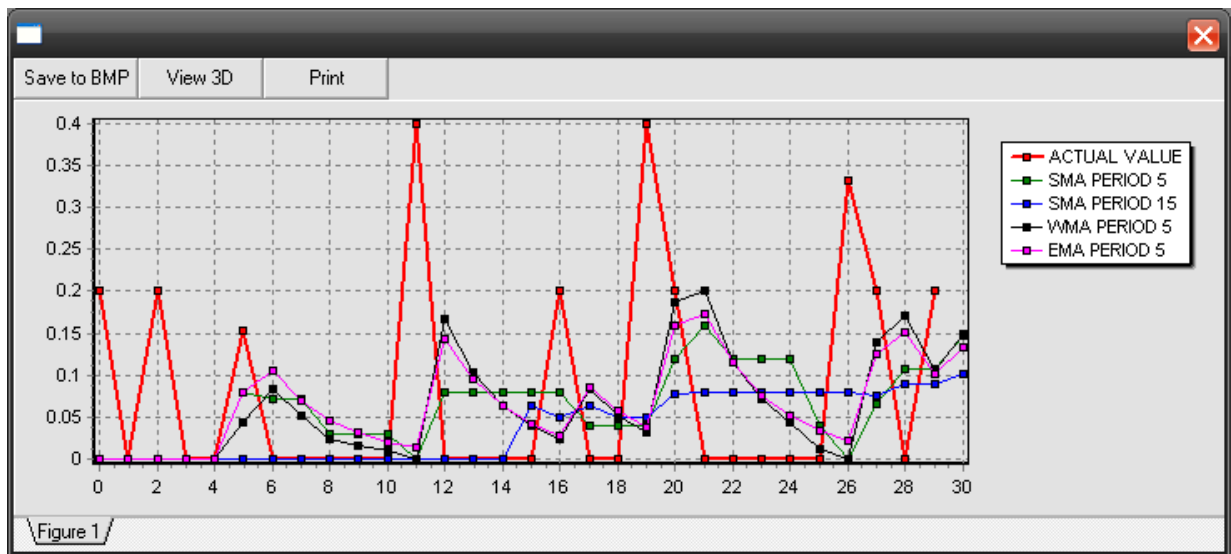
Fig. 5. The window for visualization of three prediction methods together with the base line for missing knowledge

Fig. 3, 4, and 5 illustrate the specific feature of the *SMA*, e.g. absence of several initial prediction values which number is equal to the number of time periods for calculation of the *MA* (for the first and second methods their number is 5 and 15 respectively). In all considered cases the *SMA* with 5 periods works so well as the much more difficult for calculation *EMA* with the same, e.g. 5 number of periods. This finding also is in line with the theoretical basis of prediction. The only case where the corresponding *MA* considerably diverges from each other is when the weight coefficients, assigned to the latest data, are different. The question which type of the *MA* is the best choice has no correct answer. As a rule, the *EMA* is more sensitive to changes than the *SMA*, but less than the *WMA*. However, the answer depends on the specifics of the base line, as well as on the prediction error.

**Error and Skill Analysis**

Prediction error $E_t$ for the $t$th calculation period is the difference between the actual value $Y_t$ and the predicted one $F_t$ e.g. $E_t = \left| Y_t - F_t \right|$. The evaluation of the prediction can be equally performed through analyzing certain measures of the aggregate error that could be one of the following:

➢ *Mean Absolute Error (MAE)* is the average of the difference between predicted and actual values in all test cases, e.g. it is the average prediction error and is calculated as $MAE = \dfrac{1}{N} \sum\limits_{t=1}^{N} | E_t |$.

➢ *Mean Absolute Percentage Error* (*MAPE*) is the mean or average of the absolute percentage errors of predictions as it is calculated as $MAPE = \dfrac{1}{N} \sum\limits_{t=1}^{N} | E_t / Y_t |$.

➢ *Mean Absolute Deviation Percentage* (*MADP*) is calculated refer to the proportion

$$MADP = (\sum_{t=1}^{N} |E_t|) / \sum_{t=1}^{N} |Y_t|.$$

➢ *Mean Squared Error* (*MSE*) is the average loss, e.g. the expectation of the squared deviations of the arguments from their respective target value. It is calculated as $MSE = \dfrac{1}{N} \sum_{t=1}^{N} E_t^{2}$ .

➢ *Root Mean Squared Error* (*RMSE*) is one of the most commonly used measures of success for numerical prediction and is computed by taking the average of the squared differences between each predicted value and its corresponding correct value, e.g. $RMSE = \sqrt{MSE} = \sqrt{(\sum_{t=1}^{N} E_t^{2}) / N}$ .

➢ *Skill in Prediction* (*SP*) is defined a root mean squared error as scaled representation of prediction error that relates the prediction accuracy of a particular prediction model to some reference one. If $\overline{Y}$ is the prediction for the period $t$ then the *SP* is calculated as:

$$SP = 1 - \frac{MSE_f}{MSE_c} ; \ MSE_f = \frac{1}{N} \sum_{t=1}^{N} E_t^{2} ; \ MSE_c = \frac{1}{N} \sum_{t=1}^{N} (\overline{Y} - Y_t)^2 ; \ \overline{Y} = \frac{1}{N} \sum_{t=1}^{N} Y_t$$

From the formula of the MAE type of error follows that the prediction is perfect if it is equal to 0. Obviously, this error will decrease with the MA decreasing. The calculated MAE error for the fourth methods in case of correct knowledge was: 0.229, 0.279, 0.228, and 0.222 respectively. In percentage that approximately means 23% that is unacceptable having in mind that the length of the scoring intervals for the marks different from "2" is 15%. In case of missing knowledge the precise values of the *MAE* were: 0.178, 0.221, 0.189, and 0.1787 respectively. In percentage that means approximately 19% that is closer to the same error of the SMA with 5 periods than to that one of the SMA with 15 periods of time. The calculated results of the MAE error in case of wrong knowledge are: 0.116, 0.119, 0.115, and 0.115 respectively. In percentage that approximately means 12% 15 % that is approximately two times lower than in case of correct knowledge. The graphical interpretation for the MAE analysis is given on fig. 6.
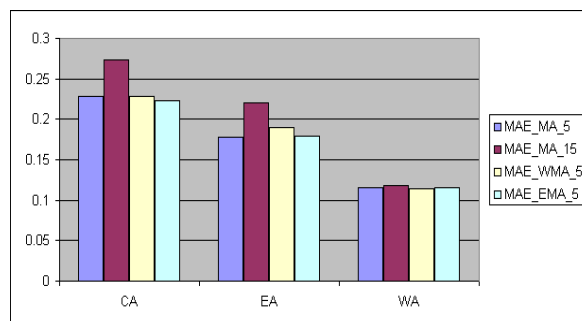


*Fig. 6.  Graphical comparison of the* MAE for all  cases

From the formula for the *MAPE* follows that this type of error can't be calculated when the base line contains a zero value. According to the formula for SP a perfect prediction has a SP equal to 1.0, a prediction with similar skill to the reference prediction would be equal to 0.0, and a prediction which is less skillful than the reference prediction will have negative values. The picture on fig. 6 is slightly changed when the most precise indicator, e.g. the *SP* is used (fig. 7). In case of correct knowledge for all four methods its value is negative and very close to zero (-0.1, -0.14, -.011, and - 0.62 respectively). In practice that means none of the methods is acceptable. In case of the missing knowledge the values of *SP* are positive (0.84, 0.69, 0.83, and 0.88 respectively) and close to a "good" prediction. The prediction is "excellent" only in case of wrong knowledge, as the values of *SP* (0.95, 0.94, 0.95, and 0.85 respectively) are positive and very close to 1.00.
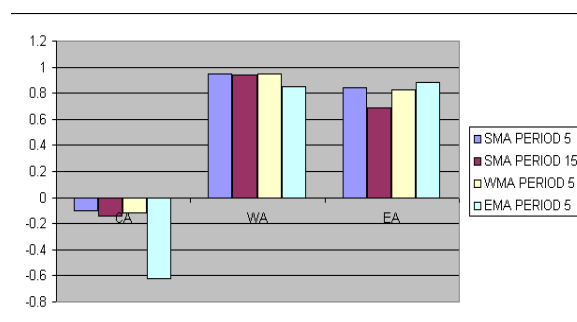


*Fig. 7. Graphical comparison of the SP for all cases*

The actual results can significantly differ from the predicted ones, which can be due to some side effects and/or external factors not taken into consideration. Regarding learners, they could be: attempt for fraud, unfamiliar type of tasks, insufficient attention or low motivation. Other external factors, related to the teacher, could be: poor session planning, organization, and/or delivery.

## Conclusions

Application of four simple for computation methods of prediction, e.g. simple moving average with 5 and 15 periods of time, weighted moving average, and exponential moving average has been illustrated for the short-term task for prediction of correct, missing, and wrong knowledge of testing.

The findings from this study are more likely not to be valid for other individual students as the process of testing as the process of learning depends to great degree on the L's attitudes. In connection with this the following methodology for the test performance prediction for other Ls is recommended using the same tool for data mining: 1) Constructing the input table with the rows equal to the test questions, and columns to the number of predicted cases, e.g. correct, missing, and wrong knowledge; 2) Adding new columns with the normalized values of the base lines; 3) Programming formulas of prediction for all chosen methods; 4) Generation of a table with the *MAE* error for the all methods; 5) Generation of a table with the *SP* values; 6) Drawing the first base line with the corresponding predicted lines by using command **PLOT;** 7) Repeating  step 4,5 and 6 for all chosen for comparison methods; 8) Making decision about preferable method for  prediction on this base respectively to the 15% length of the six-scale intervals.

**Bibliography**

[1]  Carlberg C., Business Analysis with Microsoft Excel, Sofia, "SoftPress", 2003.

[2]  Kir O., Zheliazkova I, Teodorov G., Educational Data Mining by Means of a Power Instructor's Tool, Proceedings of      International Conference on Entrepreneurship, Innovation, and Regional Development (ICEIRD), 2011 (Accepted).

[3]  Romero Cr., Ventura S. Educational Data Mining: A Review of the   State of the Art, IEEE Transactions on Systems, Man, and Cybernetics – Part C: Applications and Reviews 2010, No. 40/6, pp. 601-18.

[4]  Zheliazkova I. I., Andreeva M. H., An Intelligent Multimedia Environment for Knowledge Testing, E- learning and the Knowledge Society, Gent & Brussels, Belgium, 6-8 September 2004, pp. 3.13.1-3.13.24.

[5]  Zheliazkova I. I., Atanasov A. Z., An Approach to Teaching Systems for Decisions Making Support in Company Management, Bulgarian Journal "Automatics and Informatics", No. 1, 2008, pp. 62-68 (in Bulgarian).

[6]   Zheliazkova I. I.,  Kolev R. T., Task Results Processing for the Needs of Task-Oriented Design Environments, Int. J. Computers & Education, vol. 51, 2008, pp. 86-96.

[7]  http://www.amstat.org/publications/jse/v11n1/datasets.hays.html.

[8]  http://www.smetoolkit.org/smetoolkit/en/content/en/416/Demand-Forecasting.

[9]  http://en.wikipedia.org/wiki/Exponential_smoothing

[10] http://en.wikipedia.org/wiki/Forecast_skill.

**Authors' Information**

**Oktay Kir** – *PhD student, University of Rousse, Studentska street  8, Rousse 7017, Bulgaria;*

*e-mail: kir.oktay@gmail.com*

**Irina Zheliazkova** – *Associate Professor; University of Rousse, Studentska street  8, Rousse 7017, Bulgaria; e-mail: irina@ecs.ru.acad.bg*