# Data-driven English learner level management: K-means clustering algorithm improvement and application

**Wenjing Huang[1],***

[1] School of Foreign Languages, Hubei Engineering University, Xiaogan, Hubei, 432000, China

Corresponding authors: (e-mail: wenjingh91@163.com).

**Abstract** As an important means to improve the efficiency of English teaching, the current method of student stratification management still has shortcomings such as being too broad and not fully considering the students' situation. This paper organizes the relevant concepts of Bayesian network, and on the basis of this theory, proposes Bayesian network prediction as a prediction method of students' performance. At the same time, in order to analyze and predict students' performance more scientifically and objectively, the simple Bayesian classification method is introduced. Combining the traditional weighting method and informatics method to calculate the weight value, the premise of extending the conditional independence of the simple Bayesian classifier. Based on the prediction of students' performance, the layered teaching concept is used as a framework to design the English layered teaching model and the learner's layering method. The stratified management of students is realized by adopting corresponding teaching methods based on students' English achievement characteristics. The tiered management method was recognized by 90.00% and above of the students in the practical application.

**Index Terms** layered teaching concept, English teaching, Bayesian network, achievement prediction

## I. Introduction

Nowadays, the rapid development of social economy and the deepening of the reform process of education system provide an important opportunity and a good social environment for the improvement of the quality and efficiency of college English teaching. As an important content of the new curriculum reform model under quality education, the college classroom has put forward higher classroom teaching requirements for college English [1]. On the basis of the innovative English teaching mode, not only should we focus on the classroom education and teaching methods, but we should also take students as the basis, combine the actual development of students, and carry out layered teaching [2]-[4]. Stratified teaching is a new attempt under the educational concept of the new curriculum reform. The teaching mode embodies the teaching idea of modern education and student development-oriented, and further interprets the teaching concept of teaching according to students' abilities, and gradually explores the teaching methods suitable for students' actual situation in the teaching process, so as to continuously improve the efficiency of college English teaching [5]-[7].

English as a two-way teaching activities, the need for teachers and students of communication linkage, to build a new type of teacher-student relationship to promote the development of students [8]. Teachers stratify the lessons according to the students' knowledge understanding in the English classroom and develop instructive and diversified teaching strategies by combining the students' interests as well as their personality traits, which will effectively strengthen and improve the deficiencies existing in traditional English education [9], [10]. At the same time, teaching measures based on different students' cognitive levels and their own learning abilities are established to give full play to students' strengths and make up for their weaknesses, so that students can maximize their own improvement in the process of recent development [11], [12]. In the new era, English teachers can use intelligent algorithms to manage students hierarchically, which reflects the student-oriented teaching concept, and by dividing students according to the actual situation, different educational strategies can be carried out for different levels of students [13]-[15]. It can be seen that the construction of English layered teaching from the perspective of student management is crucial to the improvement of teaching efficiency and the cultivation of students' comprehensive ability, and it is a necessary way to cater for quality education.

This paper firstly elaborates the relevant conceptual content of Bayesian network theory in detail, and then proposes an English performance prediction method based on Bayesian network. Secondly, we choose the plain Bayesian classification method as the classification method of students' performance, integrate the traditional weighting method with the informatics method, determine the comprehensive weights, and build the plain Bayesian

classifier. At the same time, we design experiments to compare the predicted grades with the real grades and the performance comparison with similar methods to test the effectiveness of the grade prediction method based on Bayesian classification algorithm. Based on the prediction results of students' performance, we design a hierarchical teaching process and a learner hierarchical method to carry out student hierarchical management in English teaching. Finally, the designed hierarchical management method is applied to actual English teaching to evaluate the feasibility and reliability of the method from the perspectives of students' performance, changes in learning problem solving, and usefulness to students.

## II. Grade prediction based on Bayesian classification algorithm

### II. A. Concepts related to Bayesian networks

Full probability formula: Let $\{B_1, B_2, \cdots, B_n\}$ ($n$ is either finite or infinite) be a cluster of complete events (also known as a division of $\Omega$) in the sample space, in other words, they must satisfy the following conditions:

(1) They are mutually disjoint, i.e., $B_i B_j = \phi$, $i \neq j$.

(2) Their sum (sum) is exactly the sample space, i.e. $\sum_{i=1}^{n} B_i = \Omega$. Let $A$ be an individual event in $\Omega$, then the full probability formula is equation (1):

$$P(A) = P\left(\sum_{i=1}^{n} AB_i\right) = \sum_{i=1}^{n} P(A \mid B_i)P(B_i) \tag{1}$$

The probability of occurrence of various events based on information from past history or subjective will judgment is called a priori probability.

Bayes formula (Bayes formula): under the condition of the full probability formula, i.e., there exists a complete event group $\{B_1, B_2, \cdots, B_n\}$ in the sample space $\Omega$, let $A$ be an event in $\Omega$ and $P(B_i) > 0$, $i = 1, 2, \cdots, n$, and $P(A) > 0$, then there is Equation (2) according to the conditional probability calculation method:

$$P(B_i \mid A) = \frac{P(A \mid B_i)P(B_i)}{\sum_{j=1}^{n} P(A \mid B_j)P(B_j)} \tag{2}$$

Then using the Bayesian formula, new information is obtained in conjunction with the survey analysis and the more realistic probabilities obtained by correcting the a priori probabilities are called posterior probabilities.

The centerpiece in Bayesian statistical methods is the Bayes formula. In addition to the need to utilize sample information, Bayesian statistical methods must also utilize prior information about the parameters, and when there is enough prior information, a prior distribution can be obtained. The prior distribution $\pi(\theta)$ is the knowledge of the possible values of the parameter $\theta$ before the sample $X$ is drawn. After obtaining the sample data, since the sample $X$ also contains information about $\theta$, so when once the sampling information data, the knowledge of $\theta$ occurs to be changed and adjusted to it, the result of the adjustment to obtain a new knowledge of $\theta$, known as the posterior distribution, notated as $\pi(\theta \mid x)$. Therefore, the posterior distribution can usually be viewed as the result of the adjustments one makes to the prior distribution using both aggregate and present information (also collectively referred to as sampling information). In other words, the posterior score is a synthesis of the three types of information.

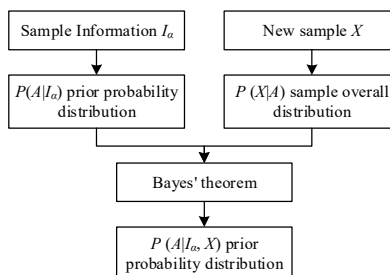The Bayesian statistical process is illustrated in Figure 1.



Figure 1: Bayesian statistical process

## II. B.Bayesian network prediction methods

Bayesian prediction specifies a prior probability distribution for the unknown parameters of the overall distribution in the prediction process, also as we gave the definition earlier, so that the posterior probability distribution (predicting the distribution of student achievement levels) can be computed by using the prior probability distribution, the overall distribution (distribution of student achievement levels), the sample information (statistical data related to the students) and combining it with Bayesian formulas, which are the the prediction target we need. The general form of the Bayesian prediction model is shown in Figure 2.
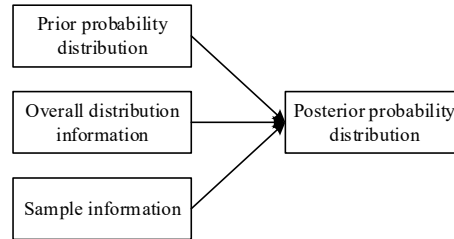


Figure 2: Bayesian prediction idea

Another idea of Bayesian statistical methods for prediction is to build dynamic models and to view the predictive distribution as a conditional probability distribution, which is what is used in this paper. We can base on the prior distribution information $\xi_t \mid D_{t-1}$, in this paper, the prior distribution refers to the distribution of the grade of the student's performance, and combined with the Bayesian formula to find the predictive distribution $p(y_t \mid D_{t-1})$, and the paper needs to get the predictive distribution is the distribution of students' academic performance under the ten impact factors. The paper needs to get the predictive distribution of the students' academic achievement level under the ten impact factors, and also apply the Bayesian formula to calculate the posterior distribution information $\xi_t \mid D_t$, and after that keep on correcting the prior information to get the predicted values. The framework of which Bayesian prediction recursive algorithm is shown in Fig. 3.
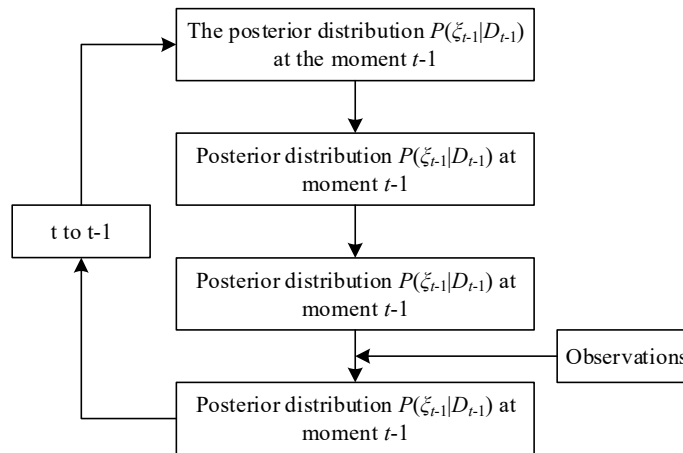


Figure 3: Bayesian predictive recursive algorithm

The Bayesian network model is a mathematical model based on probabilistic reasoning, which uses a graphical network to intuitively represent the joint probability and conditional independence of variables, which is very useful for probabilistic reasoning. The Bayesian network model is a directed acyclic graph (DAG) with $N$ node, the points in the graph represent random variable $X = \{x_1, x_2, \cdots, x_n\}$, the directed edges between points represent nodes $X_i$ conditionally independent of the subset of points composed of non-$X_i$ descendant points of the parent point of $X_i$, $\Pi_i$ represents the parent set of variable $X_i$, $\pi_i$ represents the configuration of $\Pi_i$, $Pa_i$ represents the relationship between points, and conditional probability expresses the correlation between nodes and its parent node, which is represented by $P(x_i \mid \pi_i)$. A Bayesian network specifies any one-to-one specific

configuration in the graph. For each $X_i$ there will be a subset $\Pi_i \subseteq \{X_1, \cdots, X_{i-1}\}$ such that $X_i$ is the same as $A(X_i) = \{X_1, \cdots, X_{i-1}\} \setminus \Pi_i$ is conditionally independent given $\Pi_i$. Then, for any $X$ there will be Eq. (3):

$$P(X_i \mid X_1, \cdots, X_{i-1}) = P(X_i \mid \pi_i) \tag{3}$$

Therefore, there is equation (4):

$$P(X_i) = \prod_{i=1}^{n} p(X_i \mid \pi_i) = \prod_{i=1}^{n} P(X_i \mid X_1, \cdots, X_{i-1}) \tag{4}$$

It is clear that any joint probability distribution can be inferred from a sufficient number of conditional probabilities.

### II. C. Plain Bayesian Classifier
#### II. C. 1)    Plain Bayesian classification methods
Suppose that $D$ is a collection of training tuples and associated class labels, and that each tuple is described by a $n$-dimensional attribute vector $X = \{x_1, x_2, \cdots, x_n\}$ describing the $n$-dimensional measurement of the tuple's $n$ attributes $A_1, A_2, \cdots, A_n$ on the tuple of $n$ measurements. Assuming that there are $m$ classes $C_1, C_2, \cdots, C_m$, there is a computational equation (5) for a given tuple $X$:

$$P(C_i \mid X) = \frac{P(X \mid C_i) P(C_i)}{P(X)} \tag{5}$$

Combined with the assumption of conditional independence of classes, where $P(X \mid C_i) = \prod_{k=1}^{n} P(x_k \mid C_i)$. The classification will predict that $X$ belongs to the class with the highest a posteriori probability (under condition $X$), that is, the Plain Bayesian classification predicts that $X$ belongs to the class $C_i$ if and only if $P(C_i \mid X) > P(C_j \mid X)$, $1 \le j \le m$, $j \ne i$.

#### II. C. 2)    Weighted Bayesian Classification Methods and Weight Determination
Plain Bayesian classification has a prerequisite: classes have conditional independence, i.e., all conditional attributes have the same classification importance (weight of 1) on the decision attribute. However, this is not the case in practice; some conditional attributes have a large impact on classification and some have a smaller one. Therefore, it is necessary to pay different weights to the attributes according to their importance, and the weighted plain Bayesian classification method (WNB) is shown in equation (6):

$$V_{wnb}(X) = \arg \max_{C} P(C) \prod_{k=1}^{n} P(X_k \mid C)^{w_k} \tag{6}$$

where $w_k$ is the weight of the $k$th attribute $A_k$. The greater the weight, the greater the influence of the attribute on the classification.

For weighted Bayesian classification, the most important issue is the determination of the weights.

(1) Weights in the traditional view

Assuming that $SGF(\alpha_i, D)$ represents the importance of attribute $A_i$ relative to decision $D$, and there are a total of $n$ attributes, the weight value of $A_i$ is Equation (7):

$$w_{1i} = \frac{SGF(\alpha_i, D)}{\dfrac{1}{n} \sum_{k=1}^{n} SGF(\alpha_i, D)} \tag{7}$$

(2) Attribute weights in the informatics viewpoint

Assuming that $I(\alpha_i, D)$ denotes the mutual information of the conditional attribute $A_i$ and the decision $D$, the weights of $A_i$ are given in equation (8):

$$w_{2i} = \frac{I(\alpha_i, D)}{\frac{1}{n}\sum_{k=1}^{n} I(\alpha_i, D)} \tag{8}$$

(3) Combined weights determination

Since the traditional definition of attribute importance is complementary to the information definition, the weights are defined as the average of the weights under the above two views, i.e., equation (9):

$$w_i = \frac{w_{1i} + w_{2i}}{2} \tag{9}$$

## II. D.Experimental validation and analysis

In order to verify the validity and accuracy of the prediction method (NB) proposed in the previous paper, 25 students were randomly selected from the test dataset to predict their final grades and compared and analyzed with the real grades. Figure 4 shows the situation between the students' predicted grades and the real grades. According to the experimental results, the grade prediction method proposed in this paper shows high precision and accuracy, and the predicted final grades are relatively more stable with smaller errors, so the method is suitable for the current grade prediction needs.
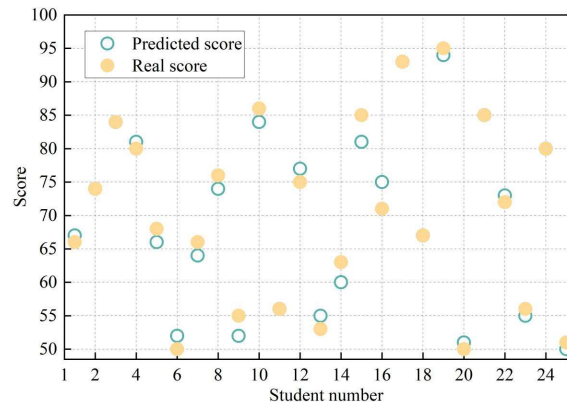


Figure 4: Performance prediction result

Under other conditions unchanged, this paper's method (NB) is compared and analyzed experimentally with three traditional prediction methods: decision tree (LR), BP neural network (BP), and support vector machine (SVM), and the comparisons of the four methods' accuracy (ACC), precision (Pre), recall (Rec), and F1 values are shown in Fig. 5, in which the performance of the four indexes of this paper's method are all greater than or equal to 0.80, far more than the other three similar algorithms, and relatively stable. It can be seen that the method in this paper has higher precision and accuracy than the traditional method, and the predicted final grades are not only more stable but also with less error, which is suitable for the prediction needs at this stage.
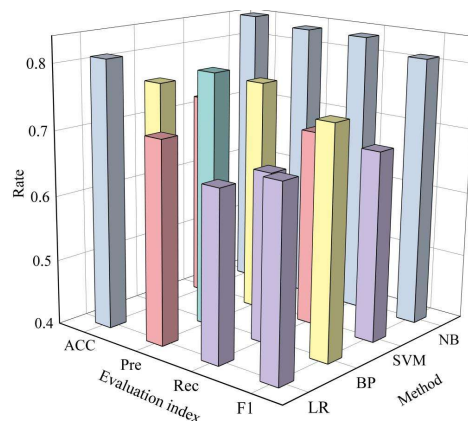


Figure 5: Comparison of algorithm performance indicators

In order to ensure the reliability of the experimental results assessment, a five-fold cross-validation method is used for further testing. Specifically, the dataset is evenly divided into five parts, one of which is selected as the validation set each time, and the remaining four parts are used as the training set, where different machine learning algorithms are used to train out the model and test the model effect on the validation set. This process will be looped five times, and each time a prediction will be obtained on the test set, and finally the five results will be averaged to get the final model prediction. By doing this, the student feature data can be collected five times and five different training and validation datasets can be obtained, thus ensuring the stability of the assessment results. The results of the five-fold cross-validation experiments are shown in Fig. 6. Compared with the other three traditional prediction methods, the prediction method proposed in this paper has a significant reduction in the MAE error index, which ranges from 4.15-4.25, which significantly improves the overall prediction effect.
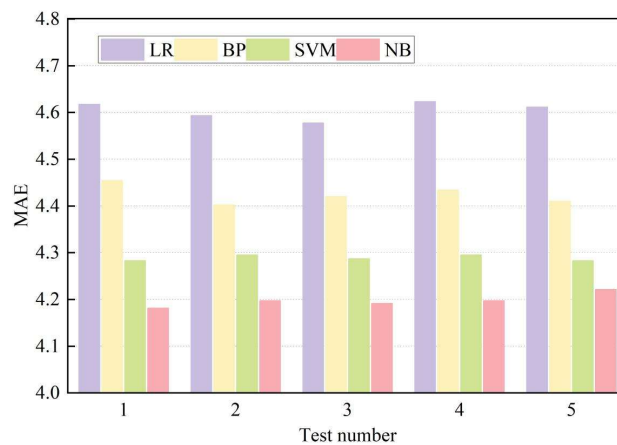


Figure 6: Comparison of MAE results

## III. Design and Application of Tiered Instruction in English Courses

In this paper, based on the current English course performance data and student performance characteristics, we predict the classification results of the target course performance and implement tiered teaching for the classification results. The application of the tiered teaching method based on the predicted results in the design process of higher education courses is shown in Fig. 7, which is based on the concept of tiered teaching, to understand the design ideas of the tiered teaching model, and to modularize the design of each tiered module.
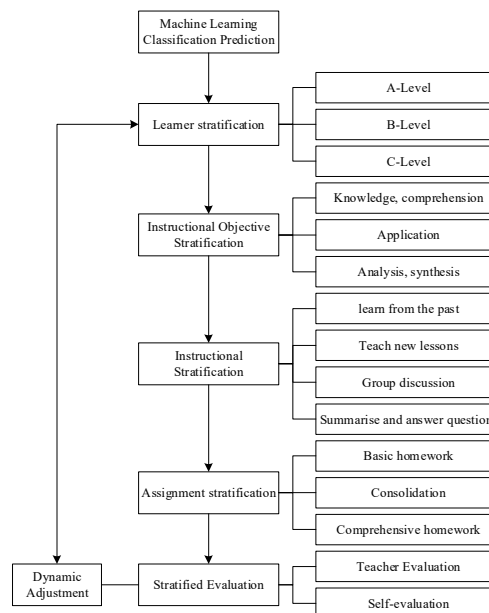


Figure 7: The hierarchical teaching method process based on the prediction results

### III. A. Stratified learners

The machine learning model is trained by the course grades obtained by students of the same major in the previous level, and the sample students' grades of the previous semester's courses similar to English application are inputted into the trained model for grade prediction, and finally the students who may obtain 80 points and above are classified as level A, those who may obtain 60-80 points are classified as level B, and those who may fail are classified as level C. The machine learning model can also be used to predict the grades of students in the same major in the previous semester. In addition, for each level of students can also be divided into a number of groups, each group by the students spontaneously elected a group leader, to assume the organization of group discussions, collect group members of homework and maintain group discipline and other responsibilities.

The stratification of students is not to differentiate students into "excellent", "good", "poor" three levels, but to different levels of students for targeted counseling, and in effective the purpose is to provide targeted counseling to students at different levels in order to effectively improve their performance while ensuring that their self-esteem is not jeopardized.

### III. B. Application Analysis and Effectiveness Evaluation

#### III. B. 1) Comparison of performance before and after experimentation

Attitude towards completing English homework, active learning of English outside class, and perceived changes in the difficulty of English learning were compared between the performance of the pre-test experimental class (BE), the performance of the post-test experimental class (AE), and the performance of the pre-test control class (BC), and the performance of the post-test control class (AC), respectively.

The changes in the attitudes of students in the experimental and control classes in completing English homework before and after the experiment are shown in Fig. 8, setting the attitudes towards completing English homework as (a1) completely perfunctory, (a2) barely coping with it, (a3) completing it conscientiously, (a4) trying their best to do it, and (a5) striving for excellence. The attitude of the experimental class towards English homework has undergone a relatively obvious and benign change. Compared with the preexperimental period, the proportion of students in the experimental class who chose (a4) try their best to do the best they can has increased by 17.28% after the experiment, while on the contrary, the change in the control class is relatively small.
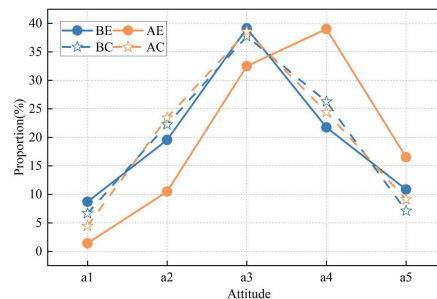


Figure 8: The change in attitude towards completing English homework

The changes in the situation of students' active learning of English outside class in the experimental and control classes before and after the experiment are shown in Fig. 9, and the attitudes that set the situation of active learning of English outside class are (b1) never, (b2) seldom, (b3) sometimes, (b4) often, and (b5) always. Compared with before the experiment, 32.45% of the students in the experimental class after the experiment often take the initiative to learn English outside the classroom, are able to be more active in searching for and learning new English knowledge or skills outside the classroom, and their initiative and self-consciousness in learning are greatly enhanced.
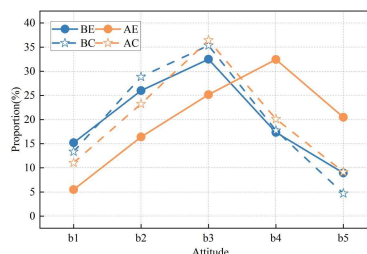


Figure 9: The changes in the situation of active English learning outside of class

The changes in the perception of English learning difficulty between the experimental class and the control class before and after the experiment are shown in Fig. 10, and the perception of English learning difficulty in the setting class are (c1) very difficult, (c2) difficult, (c3) moderate, (c4) easy, and (c5) very easy. Compared with the control class, the students in the experimental class were much less intimidated by English learning. After the experiment, only 9.56% of the students think it is very difficult to learn English. It can be seen that the method of this paper can effectively realize teaching according to students' abilities, thus reducing students' learning difficulties.
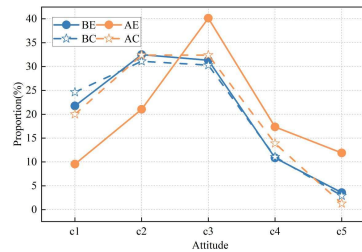


Figure 10: Cognitive changes in the difficulty of learning English

### III. B. 2) Comparative analysis of learning problem solving

The solution of learning problems is a direct manifestation of verifying the validity of the student stratified management method based on plain Bayesian classification. Two classes in a middle school were selected as experimental subjects, each with 40 students, one of which was an experimental class and was taught using the student hierarchical management method. The other class is the control class, which adopts the general teaching mode. The students of both experimental and control classes were divided into eight groups, and four of them were chosen to choose the assessment, set up with 20 problems, and tested in the third week and the seventh week respectively. The problem solving of the experimental class is shown in Table 1, and the problem ending of the control class is shown in Table 2.

Table 1: The problem-solving situation of the experimental class

| Weeks | Group | The number of problems solved | Average | Proportion |
|---|---|---|---|---|
| 3 | 1 | 16 | 16.5 | 80.00 |
| | 2 | 18 | | 90.00 |
| | 3 | 15 | | 75.00 |
| | 4 | 17 | | 85.00 |
| 6 | 1 | 17 | 17.75 | 85.00 |
| | 2 | 19 | | 95.00 |
| | 3 | 17 | | 85.00 |
| | 4 | 18 | | 90.00 |

Table 2: The problem-solving situation of the control class

| Weeks | Group | The number of problems solved | Average | Proportion |
|---|---|---|---|---|
| 3 | 1 | 13 | 12.5 | 65.00 |
| | 2 | 14 | | 70.00 |
| | 3 | 11 | | 55.00 |
| | 4 | 12 | | 60.00 |
| 6 | 1 | 10 | 10.5 | 50.00 |
| | 2 | 11 | | 55.00 |
| | 3 | 12 | | 60.00 |
| | 4 | 9 | | 45.00 |

Comparing Tables 1 and 2, the mean value of problem solving in the experimental class is significantly higher than that of the control class and shows an upward trend, with the number of problems solved amounting to 17.75 at week 6. While the mean value of the control class is lower and shows a downward trend, decreasing from 12.5 problems at week 3 to 10.5 problems at week 6. It indicates that the student hierarchical management method

based on plain Bayesian classification is going to be suitable for the students and more helpful in learning problem solving.

### III. B. 3)  Practicality survey for students

The heterogeneous grouping method of collaborative learning based on plain Bayesian classification assigns categories to each student so that each student has his/her own category in the collaborative learning group, this survey is mainly to determine whether the students can play the role of the category in the group after they have their own category, and the applicability of the student's category in the group. Therefore, this satisfaction survey set three questions on the applicability of students' categories as follows:

Q1: The plain Bayesian categorization method can help me better identify my position in the group.

Q2: Plain Bayesian categorization makes me more united with others.

Q3: The categories categorized by the Plain Bayesian classification algorithm are well suited to myself.

The chosen indicators of usefulness to the students' categories are (I1) very satisfied (%), (I2) satisfied (%), (I3) average (%), (I4) dissatisfied (%), (I5) very dissatisfied (%). The results of the survey for the 3 questions are shown in Table 3.

Table 3: Practicality survey results

| Question | I1 | I2 | I3 | I4 | I5 |
|----------|------|------|------|------|------|
| Q1 | 32.75 | 53.44 | 11.48 | 2.2 | 0.13 |
| Q2 | 24.13 | 58.06 | 10.8 | 5.18 | 1.83 |
| Q3 | 15.51 | 67.23 | 13.78 | 3.44 | 0.04 |

By comparing the data in the table, it can be seen that most of the students think that the plain Bayesian classification can help him better identify his position in the group, and only a very small number of people do not think so, and for the article that the categories classified by the plain Bayesian classification algorithm are well suited to myself, most of the people expressed their satisfaction, and a few of them expressed their dissatisfaction, and analyzing the results of the survey of the satisfaction of the two articles, it is probably due to the fact that the plain Bayesian classification is not suitable for individual students in the selection of the training dataset, there is a certain error in constructing the plain Bayes classifier, which makes the classified categories have bias in individual students. From the analysis of the results of the applicability survey of the categories classified by Park Bayes, 97.67% of the students subjectively believed that the categories classified by Park Bayes helped them to identify their position in the group, which was very suitable for themselves, and 92.99% of the students subjectively believed that the Park Bayesian classification made them more united with others. The synthesis of the survey results can be seen that the classification made by Parsimonious Bayes has a small error and is basically suitable for most of the students.

## IV.  Conclusion

This paper predicts and classifies students' English performance by combining Bayesian network prediction method and weighted Bayesian classification method. On this basis, it proposes a tiered teaching method for English courses, and then develops scientific student tier management. The grade prediction method based on the Bayesian classification algorithm has an accuracy, precision, recall, and F1 value of more than 0.80 compared with similar methods, and it combines high-precision prediction ability and stability. The designed student hierarchical management method, when applied in practice, improved the attitudes of students in the experimental class in completing English homework by 17.28%, and promoted 32.45% of the students in the experimental class to study English actively outside the classroom. The students' problem solving ability also increased from an average of 16.5 problems to an average of 17.75 problems, which was recognized by 90.00% of the students and above.

## References

[1]  Arianto, R. S., Juhana, J., & Ruminda, R. (2023). Building students' confidence in speaking English through Differentiated Instruction. Lectura: Jurnal Pendidikan, 14(2), 276-287.

[2]  Ismail, S. A. A., & Al Allaq, K. (2019). The nature of cooperative learning and differentiated instruction practices in English classes. Sage Open, 9(2), 2158244019856450.

[3]  Sapan, M., & Mede, E. (2022). The effects of differentiated instruction (DI) on achievement, motivation, and autonomy among English learners. Iranian Journal of Language Teaching Research, 10(1), 127-144.

[4]  Suryati, I., Ratih, K., & Maryadi, M. (2023). Teachers' challenges in implementing differentiated instruction in teaching English at junior high school. Eduvest-Journal of Universal Studies, 3(9), 1693-1708.

[5]     Raza, K. (2020). Differentiated instruction in English language teaching: Insights into the implementation of Raza's teaching adaptation model in Canadian ESL classrooms. TESL Ontario Contact Magazine, 46(2), 41-50.

[6]     Hoang, T. H., & Pham, T. N. T. (2018). Hierarchy in high school English classrooms in Vietnam. Identity, equity and social justice in Asian Pacific Education. Monash University Publishing.

[7]     He, B. (2020, April). Empirical Research on College English Teaching Mode of Hierarchical and Three-Dimensional Interaction. In 5th International Conference on Social Sciences and Economic Development (ICSSED 2020) (pp. 131-134). Atlantis Press.

[8]     Anggraeny, T. F., & Dewi, D. N. (2023). Analysis of teacher strategies in teaching English using differentiated learning. EJI (English Journal of Indragiri): Studies in Education, Literature, and Linguistics, 7(1), 129-146.

[9]     Kupchyk, L., & Litvinchuk, A. (2020). Differentiated instruction in English learning, teaching and assessment in non-language universities. Advanced Education, 89-96.

[10]    Kamarulzaman, M., Azman, H., & Zahidi, A. (2017). Differentiated instruction strategies in English language teaching for gifted students. Journal of Applied Environmental and Biological Sciences, 7(1), 78-90.

[11]    Grecu, Y. V. (2023). Differentiated instruction: Curriculum and resources provide a roadmap to help English teachers meet students' needs. Teaching and Teacher Education, 125, 104064.

[12]    Kamarulzaman, M. H., Azman, H., & Zahidi, A. M. (2015). Differentiation practices among the English teachers at PERMATApintar national gifted and talented center. Asian social science, 11(9), 346.

[13]    Zhang, R., & Huang, H. (2022). Application of Convolutional Neural Network‐Based Hierarchical Teaching Method in College English Teaching and Examination Reform. Mathematical Problems in Engineering, 2022(1), 3378599.

[14]    Zhang, C. (2022, July). The application of hierarchical teaching mode based on hybrid criterion fuzzy algorithm in higher vocational English education. In EAI International Conference, BigIoT-EDU (pp. 424-430). Cham: Springer Nature Switzerland.

[15]    Yao, R. (2021, March). Construction of bilingual teaching hierarchy model based on flipping classroom concept. In International Conference on Data and Information in Online (pp. 421-425). Cham: Springer International Publishing.