

A novel model to predict mental distress among medical graduate students in China.

Fei Guo

The Second Affiliated Hospital of Air Force Medical University

Min Yi

The First Affiliated Hospital of the Southwest Medical University

Li Sun

The Second Affiliated Hospital of Air Force Medical University

Ting Luo

The Second Xiangya Hospital of Central South University

Ruili Han

The Second Affiliated Hospital of Air Force Medical University

Lanlan Zheng

The Second Affiliated Hospital of Air Force Medical University

Shengyang Jin

Chinese Academy of Medical Sciences and Peking Union Medical College

Jun Wang

Shaanxi Provincial Armed Police Hospital

Mingxing Lei

Hainan Hospital of Chinese PLA General Hospital

Changjun Gao (✉ gaocj74@163.com)

The Second Affiliated Hospital of Air Force Medical University

Research Article

Keywords: Mental distress, Prediction model, Medical graduate student.

Posted Date: April 12th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-394947/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License. [Read Full License](#)

Version of Record: A version of this preprint was published at BMC Psychiatry on November 15th, 2021. See the published version at <https://doi.org/10.1186/s12888-021-03573-9>.

Abstract

Background: Several studies have reported serious mental status among medical graduate students, which triggered a negative impact on their physical and psychological health. This study aimed to develop a novel prediction model to calculate the risk of mental distress among medical graduate students.

Methods: This study analyzed 1079 graduate students via an online questionnaire. Included subjects were randomly divided into the training group and validation group. In the training group, a formula was developed, and validation of the formula was performed in the validation group. The discrimination and calibration ability were assessed for the predictive performance of the formula.

Results: One thousand and fifteen subjects were enrolled and randomly divided into the training group (n=508) and the validation group (n=507). The prevalence of severe mental distress was 14.96% in the training group, and 16.77% in the validation group. The formula included six variables, including year of study, type of student, daily research time, monthly income, scientific learning style, and feeling of time stress. The area under the receiver operating characteristic curve (AUROC) and calibration slope for the formula were 0.70 and 0.90 (95% CI: 0.65~1.15) in the training group, respectively; and 0.66 and 0.80 (95% CI: 0.51~1.09) in the validation group, respectively.

Conclusion: Six risk factors for anxiety and depression were identified and a prediction model was created. The formula may be a useful model that can identify a high risk of mental distress among medical students.

Background

Mental distress, characterized by a broad range of behavioral and psycho-physiological symptoms, is a mental health problem often relating to mental disorders such as depression and anxiety [1]. A global prevalence of depression among medical students was up to 28.0% [2]. The prevalence of anxiety was 33.8% among medical students globally, which is significantly higher than the general population [3]. Moreover, a recent study reported high prevalences of depression (29%) and anxiety (21%) among Chinese medical students [4]. All these data alarmed that a large number of medical students around the world were experiencing severe mental distress which could impair their psychosocial functioning, physical health, professional and academic performance, and ultimately cause serious consequences including divorce, crime, self-harm, and suicide tendency [5-9]. Meanwhile, a large body of literature revealed that the mental distress was the largest cost driver of the global economic burden of non-communicable diseases [10]. Therefore, it is an urgent issue to find causes and solutions of mental distress.

Screening risk factors of anxiety and depression helps early detection and intervention, and prevent more serious consequences. As is indicated in recent studies, medical students are at high risks of mental distress, which are imputed to the high levels of academical, psychological, and emotional stress including academic demands, workload, pressure from teachers and parents, financial burden and worry about the future [11-13]. These factors were concluded as medical environment, ethical conflicts, exposure to suffering, life events and educational debt, and applied in kinds of testing scales including the Symptom Check List-90 (SCL-90), Beck Depression Inventory (BDI), Center for Epidemiological Studies Depression Scale (CES-D), Hopkins Symptom Checklist (HSCL), State-Trait Anxiety Inventory (STAI), Profile of Mood States (POMS), and Brief Symptom

Inventory (BSI). However, these factors cannot be quantified, and the association with mental distress cannot be measured. In addition, the above scales are primarily indicated for professionals to evaluate the mental status of general population who are often unaware of the abnormal changes. As a former study indicated, the role of inadequate self-awareness about one's mental health concerns was a barrier to reaching out for professional help [14]. It highlights the importance of expanding the range of factors beyond commonly studied concepts like the demand-control model and the effort-reward imbalance model [15].

Therefore, this study aimed to develop a novel model to present the relative attribution of both mental distress and potential risk factors which can be detected and utilized conveniently.

Methods

Study design and sample size estimation

A cross-sectional survey was conducted from November to December 2020 among medical graduate students. An informed consent was obtained from all subjects that were asked to complete an online questionnaire voluntarily without any financial compensation. All valid information including associated device, IP address and answers for each question was collected anonymously, and then constructed a basic database about the mental distress among medical graduate students by automatic collation and graphical representation for each question.

We enrolled subjects apart from those who were post-doctors and reported a diagnosis with depression or anxiety. The subjects were then randomly divided into the training group and validation group. The training group was used to develop a formula to calculate the prevalence of mental distress among medical graduate students in China. Meanwhile, internal validation of the formula was performed in the validation group.

For the estimation of sample size, we took the prevalence of 28% [16] for mental distress from a study done among Chinese graduate students, 95% certainty and $\pm 5\%$ margin of error and using the population correction formula. Considering 10% of non-response rate, the sample size was estimated to be 344.

Ethics approval and study registration

The aims and procedures of the study were reviewed and approved by the Research Ethics Committee of Plastic Surgery Hospital of Chinese Academy of Medical Science (Approval number: 2020157). The study was registered at the Chinese Clinical Trial Registry (Registration number: ChiCTR2000039574). All procedures used complied with the ethical principles on human experimentation and with the Helsinki Declaration of 1975 as revised in 2008.

Instruments

Data on potential risk factors and the main observation (severe mental distress) in this study were collected by a questionnaire consists of sociodemographic characteristics, academic performance, incumbency of tutor, and psychological evaluation. Sociodemographic characteristics contains age, year of study, major and school location (provincial capital or other cities), marital status and monthly income. Academic performance includes degree pursuing, university, type of student, kinds of research, daily research time, scientific learning style,

number of research projects and published papers, feeling of time stress (range from 1-7, 1 means none, 2-3 mild, 4-5 moderate, and 6-7 severe). Incumbencies of tutors means tutors have positions in the department, Chinese Academy of Sciences, Chinese Academy of Engineering, or national academic organizations. Whether tutor won a bid of NSFC (National Natural Science Foundation of China) or not within the past 5 years was collected. Psychological evaluation was based on the Generalized Anxiety Disorder Scale-7 (GAD-7) and Patient Health Questionnaire-9 (PHQ-9).

GAD-7, developed by [17], is a 7-item self-report scale to measure anxiety symptoms. Each question was designed to assess the frequency of anxiety, with scores ranging from 0 (never) to 3 (daily). The total score is 0 to 21, coming from the sum of the values for each item. The reported Cronbach's α coefficient of the GAD-7 among Chinese subjects is 0.92 [18].

PHQ-9, developed by [19], is a 9-item scale based on criteria for depressive disorders in the Diagnostic and Statistical Manual of Mental Disorders (DSM-IV) [20] to measure depression symptoms. Each item scores from 0 to 3 according to increasing intensity of symptoms. The PHQ-9 had a Cronbach's α of 0.86 [21] in this study. The severe mental distress in this study was defined as the sum of GAD-7 and PHQ-9 scores ≥ 30 .

Statistical analysis

Descriptive statistics were tabulated for the overall sample and stratified by the type of answers received. Continuous variables were presented as mean \pm standard deviation (SD), while frequency and percentage were calculated for categorical variables. The potential risk factors were screened by the Least Absolute Shrinkage and Selection Operator (LASSO) method. Then, variables with a coefficient value > 0.01 were included in a multinomial logistic regression model to explore the estimates of the included variables in the formula. Statistical significance was set at $P < 0.05$ level with two-sided tests. Statistical analyses were performed using SAS 9.2 (SAS Institute Inc., Cary, NC) and R version 3.5.3 for Windows XP.

Formula development

In this study, a LASSO technique combined with the 10-fold cross-validation was used to investigate potential predictors according to computing efficient model descriptions of nonlinear systems. Variables with a coefficient value of more than 0.01 were included in the formula. The estimates used to develop the formula were obtained after the included variables re-entered the multiple logistic regression analysis. Finally, a formula was developed:

$$P(Y=1) = \frac{e^{\text{intercept}+ax_1+bx_2+\dots+ix_n}}{1+e^{\text{intercept}+ax_1+bx_2+\dots+ix_n}}$$

In the formula, a , b , ..., and i were the estimates, x_1 to x_n were the included variables, and $P(Y=1)$ indicated the predicted prevalence of severe mental distress among medical graduate students.

Validation of the formula

Internal validation of the formula was performed with the discrimination and calibration in the training and validation group. Discrimination ability of the formula was to separate students who developed mental distress from those who did not. Calibration ability of the formula was the consistency to observe and predict the prevalence of severe mental distress. The AUROC, which is the probability of concordance between predicted and observed the prevalence of mental distress among medical graduate students, was also calculated to measure the predictive effects of the formula's discrimination ability. An AUROC of more than 0.7 indicates good predictive performance and 0.8 or above indicates excellent predictive performance.

Furthermore, discrimination ability of the formula was evaluated by the discrimination slope that was defined as the difference between the mean predicted risk probability with and without mental distress among medical students. We plotted deciles of the predicted probability of severe mental distress against the observed risk of severe mental distress in each decile and fitted a smooth line. Ideally, the slope of the fitted smooth line would be close to 1 and intercepts close to 0. Besides, the Hosmer-Lemeshow goodness-of-fit test was used to evaluate the formula's calibration ability. A P-value of more than 0.05 from this test indicates good agreement between the predicted matrix and the observed matrix.

Results

A basic database was created based on 1079 valid questionnaires extracted from 1090 respondents, with 98.99% of effective response rate. Apart from 12 post-doctors and 52 students who were reported diagnoses of depression or anxiety, 1015 subjects (three times of the sample estimate) were enrolled in this study. The 1015 subjects were randomly divided into the training group (n=508) and validation group (n=507) (Figure 1).

Basic Characteristics

Of the 1015 enrolled students from the database, basic characteristics were compared between the training group (n=508) and validation group (n=507). (Table 1)

Sociodemographic issues: The mean age of these enrolled students was 25 years and more than half were single. The monthly income of nearly 60% students were 1000-3000 (Chinese Yuan, CNY) in the two groups. Most of the students were pursuing a master's degree both in the training group (75.79%) and in the validation group (76.72%). Among students in the training group, 38.39% were in the 1st year, 35.04% in the 2nd year, 21.46% in the 3rd year, and 4.92% in the 4th year or deferment period. While in the validation group, 42.41% were in the 1st year, 36.09% in the 2nd year, 19.72% in the 3rd year, and 1.78% in the 4th year or deferment period. The majority of students in the two groups were specialized in clinical fields and their school locations are out of Beijing.

Academic performance: Nearly 3/4 of the students were from the 'Double First-rate' university in both training group and validation group. A large proportion of the students in the two groups were research oriented and engaged in basic scientific research; more than one half students were working more than 6 hours per day in both training group(57.68%) and validation group(54.04%); about 70% of the students among the two groups had participated in 1-3 research projects; however, most of them (67.68%) never published an academic paper in English or in Chinese. Among the enrolled students, 56.69% in the training group were reported with a severe feeling of time stress in their scientific research, compared with 60.16% in the validation group.

Incumbency of tutor: More than 55% tutors were leaders of the department among the two groups. Only 6 tutors were academicians of the Chinese Academy of Sciences or Chinese Academy of Engineering. Twelve percent of the tutors in the training group and 8% tutors in the validation group had positions of national academic organizations. Almost two third tutors in the two groups had won the bid of National Natural Science Foundation of China in past five years.

Results of the psychological testing: According to the defined cut-offs of GAD-7 and PHQ-9, the prevalence of severe mental distress was 14.96%. The scores were 8.75 ± 5.41 for GAD-7 and 8.65 ± 6.17 for PHQ-9 in the training group. In the validation group, the prevalence of severe mental distress was 16.77%. The scores were 8.8 ± 5.39 for GAD-7 and 8.90 ± 5.94 for PHQ-9. The outcome of GAD-7 and PHQ-9 showed a good fitting in Figure 2.

The formula development

After data extraction, the 1015 subjects were randomly divided into the training group (n=508) and validation group (n=507). Seven predictors, including year of study, type of student, kinds of research, daily research time, monthly income, scientific learning style and feeling of time stress, significantly associated with the prevalence of severe mental distress, were identified by the LASSO method combined with the 10-fold cross-validation. The kinds of research was not included in the formula impute to the low coefficient value (<0.01). Finally, the left six variables were included in the formula and the corresponding estimates were obtained from the multiple logistic regression analysis (Table 2). A formula was developed: $P(Y=1) = X = -5.06 + 0.20 * \text{Year of study} + 0.44 * \text{Type of student} + 0.51 * \text{Daily research time} - 0.28 * \text{Monthly income} + 0.30 * \text{Scientific learning style} + 0.39 * \text{Feelings of time stress}$. $P(Y=1)$ indicates the predicted probability of severe mental distress. The score in each variable was assigned according to the original dataset.

Internal validation of the formula

The formula expressed a good-to-excellent discrimination ability exactly as the AUROC for itself was 0.70 in the training group and 0.66 in the validation group (Table 3 and Figure 3-4). The corresponding discrimination slopes were 0.06 (95% CI: 0.04~0.08, $P < 0.001$) and 0.04 (95% CI: 0.02~0.06, $P < 0.001$) (Figure 5-6), respectively. The correct classification rates were 82.30% in the training group, and 81.30% in the validation group. Comparing sensitivity and specificity between the training group and validation group, they were 11.80% vs. 4.70% and 94.70% vs. 96.70%, respectively. The false-positive were more than 70% and false-negative rates were both below 20% in the two groups, which indicated a high sensitivity of the formula.

The calibration slope of the formula was 0.90 (95% CI: 0.65~1.15) in the training group and 0.80 (95% CI: 0.51~1.09) in the validation group (Table 3 and Figure 7-8). Because the X- and Y- intercepts which were almost close to zero, the formula had excellent calibration ability.

Discussion

A formula was successfully developed to accurately assess the prevalence of severe mental distress depending on the basic database including 1015 subjects. The data of the 1015 subjects were used to develop and internally validate the formula which was simple and convenient since it consisted of only six variables,

including year of study, type of student, daily research time, monthly income, scientific learning style and feeling of time stress. All these variables were readily available.

After internal validated, the formula's excellent discrimination and calibration ability was confirmed according to the result (AUROC=0.66, Calibration slope of 0.04, 95% CI: 0.02~0.06, X-intercept of -0.01, 95% CI: -0.11~0.04, and Y-intercept of 0.01, 95% CI: -0.05~0.06). The P-value of Hosmer and Lemeshow Goodness-of-Fit Test was 0.97, which showed that the formula could be a reliable prediction model for the probability of mental distress among medical graduate students. From the table 4, according to the actual rate of severe mental distress of individuals, the enrolled subjects were divided into low risk (0.00%~9.99%), moderate risk (10.00%~19.99%) and high risk (20.00%~) groups, we also found a significant difference ($P \leq 0.001$) among the three risk groups that showed an excellent discrimination of the novel model.

The LASSO method combined with 10-fold cross-validation was taken to select potential risk factors in this study. Compared with other logistic regression models, it is a popular model building procedure that shrink a subset of coefficients to zero, and could perform variable selection and estimation simultaneously [22].

Over the past decades, several medical graduate students committed suicide due to mental distress each year, causing much grief and loss to the society and their families. Studies have attempted to explore normative or specific (ideographic) prediction models that are available for individuals. Allen et al. [23] developed a short-term prediction model of suicidal thoughts and behaviors, but it was applied to adolescents. Kyron et al. [24] found that short-term fluctuations in self-reported mental health may provide an indication of when an individual is at-risk of self-harm. However, the final model in this study showed marginally poorer predictive qualities when standard logistic regression was performed, with slightly lower sensitivity (71.4%), specificity (77.8%), and PPV (23.9%) statistics. Khazanov et al. [25] conducted a study on the role of distress in predicting treatment outcomes of depression, and found that assessing distress prior to treatment may help determine which patients would benefit most from adding cognitive therapy to antidepressant medications. These findings supported the generic model and the implication that it can be used as a basis to formulate and treat multiple presenting mental problems. However, assessing distress and mental risk factors may not have fully captured aspects of one's mental health [26]. As Fernandez et al. [27] stated, it was possible to develop an algorithm with good discrimination for the onset identifying overall and modifiable risks of common mental disorders among working men, but it was a secondary analysis of the study. Recently, Van Hoffen et al. included distress in a multivariable prediction model for mental long-term sickness absence (LTSA), but the external validation (an AUC of 0.70 is not sufficient) exemplified the prediction model [28]. Compared with these studies, the variables included in this study are newer, more comprehensive and more representative with excellent prediction ability, and the formula has a good fitting on the medical graduate students.

For example, for a professional oriented student (2 points) in the third year (3 points) who spent 6-10 hours (2 points) on scientific research daily with monthly income ≥ 1000 CNY (1 point), had a mild feelings of time stress (2 points), and his or her scientific learning style was guiding by others (1 point), the predicted probability of severe mental distress was

$$P(Y=1) = \frac{e^x}{1+e^x} = \frac{e^{-1.76}}{1+e^{-1.76}} = 10.00\% \quad (X = -1.76 = -5.06 + 0.20 * 3 + 0.44 * 2 + 0.51 * 2 - 0.28 * 1 + 0.30 * 1 + 0.39 * 2).$$

Limitations and implications

The formula developed in this study had both practical and theoretical implications. On one hand, it enriches and develops the findings of psychological factors linked to stress and severe mental distress in prior studies. On the other hand, it provides a new way for medical graduate students to get self-report and intervene the possible mental distress in advance.

Several limitations exist even though the current formula may be a promising prediction model of mental distress. The sample of this study was gathered by an online platform which had limited convert ability and may compromise the generalizability of the findings. The fit and relevance of the novel model should also be tested in other randomly collected samples from medical graduate students. Moreover, further evaluation and exploration of the formula should be developed by more comprehensive predictors of mental distress.

Conclusions

To our knowledge, this study is the first that develops a model to predict mental distress among medical graduate students in China. According to the results, the formula showed an excellent discrimination and calibration ability that could identify students with high risks of mental distress. Since timely screening and proper intervention was urgent among Chinese medical graduate students, this formula has the potential to be highly recommended to educational programs, mental health organizations and especially students with stigma for professional counseling.

Abbreviations

WHO: World Health Organization; GAD-7: Generalized Anxiety Disorder Scale-7; PHQ-9: Patient Health Questionnaire-9; DSM-IV: Statistical Manual of Mental Disorders; SCL-90: symptom check list-90; BDI: Beck Depression Inventory; CES-D: Center for Epidemiological Studies Depression Scale; HSCL: Hopkins Symptom Checklist; STAI: State-Trait Anxiety Inventory; LASSO: Least Absolute Shrinkage and Selection Operator; AUROC: Area under the receiver operating characteristic curve; LTSA: long-term sickness absence; CI: confident interval; CCR: correct classification rate; POS: positive; NEG: negative; CNY: Chinese Yuan; SCI: Science Citation Index; CAS: Chinese Academy of Sciences; CAE: Chinese Academy of Engineering; NSFC: National Natural Science Foundation of China.

Declarations

Acknowledgements

The authors wish to acknowledge the participation of all students in this study.

Authors' contributions

All authors have participated in the conception of the design of the study. FG, MY and LS participated in analysis and interpretation of data and drafted the manuscript. TL, RLH, LLZ, SYJ and JW participated in sending questionnaires and collected data. MXL and CJG critically reviewed the manuscript. All authors read and approved the final manuscript.

Funding

This study was supported by the National Natural Science Foundation of China (Grant No. 81971225).

Competing interests

The authors declare that they have no competing interests.

Availability of data and materials

The datasets used and/or analysed during the current study are available from the corresponding author on reasonable request.

Consent for publication

Not applicable.

Ethics approval and consent to participate

The aims and procedures of the study were reviewed and approved by the Research Ethics Committee of Plastic Surgery Hospital of Chinese Academy of Medical Science (Approval number: 2020157). An informed consent was obtained from all subjects before completing the online questionnaire.

References

1. Mirowsky J, Ross CE: **Measurement for a human science.** *J HEALTH SOC BEHAV* 2002, **43**(2):152-170.
2. Puthran R, Zhang MW, Tam WW, Ho RC: **Prevalence of depression amongst medical students: a meta-analysis.** *MED EDUC* 2016, **50**(4):456-468.
3. Quek TT, Tam WW, Tran BX, Zhang M, Zhang Z, Ho CS, Ho RC: **The Global Prevalence of Anxiety Among Medical Students: A Meta-Analysis.** *Int J Environ Res Public Health* 2019, **16**(15).
4. Zeng W, Chen R, Wang X, Zhang Q, Deng W: **Prevalence of mental health problems among medical students in China: A meta-analysis.** *Medicine (Baltimore)* 2019, **98**(18):e15337.
5. Dyrbye LN, Shanafelt TD: **Commentary: medical student distress: a call to action.** *ACAD MED* 2011, **86**(7):801-803.
6. Anckarsater H, Radovic S, Svennerlind C, Hoglund P, Radovic F: **Mental disorder is a cause of crime: the cornerstone of forensic psychiatry.** *Int J Law Psychiatry* 2009, **32**(6):342-347.
7. Twenge JM: **The age of anxiety? Birth cohort change in anxiety and neuroticism, 1952-1993.** *J PERS SOC PSYCHOL* 2000, **79**(6):1007-1021.
8. Ildstad M, Torvik FA, Borren I, Rognmo K, Roysamb E, Tambs K: **Mental distress predicts divorce over 16 years: the HUNT study.** *BMC PUBLIC HEALTH* 2015, **15**:320.
9. Tawfik DS, Profit J, Morgenthaler TI, Satele DV, Sinsky CA, Dyrbye LN, Tutty MA, West CP, Shanafelt TD: **Physician Burnout, Well-being, and Work Unit Safety Grades in Relationship to Reported Medical Errors.** *MAYO CLIN PROC* 2018, **93**(11):1571-1580.

10. Vigo D, Thornicroft G, Atun R: **Estimating the true global burden of mental illness.** *LANCET PSYCHIAT* 2016, **3**(2):171-178.
11. **Five insights from the Global Burden of Disease Study 2019.** *LANCET* 2020, **396**(10258):1135-1159.
12. Kerebih H, Ajaeb M, Hailesilassie H: **Common mental disorders among medical students in Jimma University, SouthWest Ethiopia.** *AFR HEALTH SCI* 2017, **17**(3):844-851.
13. Dyrbye LN, Thomas MR, Shanafelt TD: **Medical student distress: causes, consequences, and proposed solutions.** *MAYO CLIN PROC* 2005, **80**(12):1613-1622.
14. Dyrbye LN, Thomas MR, Shanafelt TD: **Systematic review of depression, anxiety, and other indicators of psychological distress among U.S. and Canadian medical students.** *ACAD MED* 2006, **81**(4):354-373.
15. Finne LB, Christensen JO, Knardahl S: **Psychological and social work factors as predictors of mental distress: a prospective study.** *PLOS ONE* 2014, **9**(7):e102514.
16. Guo LP, Li ZH, Chen TL, Liu GH, Fan HY, Yang KH: **The prevalence of mental distress and association with social changes among postgraduate students in China: a cross-temporal meta-analysis.** *PUBLIC HEALTH* 2020, **186**:178-184.
17. Spitzer RL, Kroenke K, Williams JB, Lowe B: **A brief measure for assessing generalized anxiety disorder: the GAD-7.** *Arch Intern Med* 2006, **166**(10):1092-1097.
18. Ying DG, Jiang S, Yang H, Zhu S: **Frequency of generalized anxiety disorder in Chinese primary care.** *POSTGRAD MED* 2010, **122**(4):32-38.
19. Kroenke K, Spitzer RL, Williams JB: **The PHQ-9: validity of a brief depression severity measure.** *J GEN INTERN MED* 2001, **16**(9):606-613.
20. Zimmerman M, Walsh E, Friedman M, Boerescu DA, Attiullah N: **Are self-report scales as effective as clinician rating scales in measuring treatment response in routine clinical practice?** *J Affect Disord* 2018, **225**:449-452.
21. Wang W, Bian Q, Zhao Y, Li X, Wang W, Du J, Zhang G, Zhou Q, Zhao M: **Reliability and validity of the Chinese version of the Patient Health Questionnaire (PHQ-9) in the general population.** *Gen Hosp Psychiatry* 2014, **36**(5):539-544.
22. Wang S, Nan B, Rosset S, Zhu J: **RANDOM LASSO.** *ANN APPL STAT* 2011, **5**(1):468-485.
23. Allen NB, Nelson BW, Brent D, Auerbach RP: **Short-term prediction of suicidal thoughts and behaviors in adolescents: Can recent developments in technology and computational science provide a breakthrough?** *J Affect Disord* 2019, **250**:163-169.
24. Kyron MJ, Hooke GR, Page AC: **Prediction and network modelling of self-harm through daily self-report and history of self-injury.** *PSYCHOL MED* 2020:1-11.
25. Khazanov GK, Xu C, Dunn BD, Cohen ZD, DeRubeis RJ, Hollon SD: **Distress and anhedonia as predictors of depression treatment outcome: A secondary analysis of a randomized clinical trial.** *BEHAV RES THER* 2020, **125**:103507.
26. Nordahl H, Odegaard IH, Hjemdal O, Wells A: **A test of the goodness of fit of the generic metacognitive model of psychopathology symptoms.** *BMC PSYCHIATRY* 2019, **19**(1):288.
27. Fernandez A, Salvador-Carulla L, Choi I, Calvo R, Harvey SB, Glozier N: **Development and validation of a prediction algorithm for the onset of common mental disorders in a working population.** *Aust N Z J Psychiatry* 2018, **52**(1):47-58.

28. van Hoffen M, Norder G, Twisk J, Roelen C: **External validation of a prediction model and decision tree for sickness absence due to mental disorders.** *Int Arch Occup Environ Health* 2020, **93**(8):1007-1012.

Tables

Table 1. Comparison of basic characteristics between the training group and validation group.			
Characteristics	Training group (508)	Validation group (507)	P
Age (mean, years)	25.56±3.09	25.32±3.26	0.22
Year of study			0.04
First year	196	215	
Second year	178	183	
Third year	109	100	
Fourth year	11	6	
In deferment period	14	3	
Major			0.28
1 Obstetrics and gynecology	29	20	
2 Surgery	96	118	
3 Internal Medicine	138	122	
4 Basic Medicine	35	32	
5 Others	210	215	
School location is in Beijing			0.68
Yes=1	51	47	
No=0	457	460	
Marital status			0.09
Single	268	245	
In love	187	184	
Married without child bearing	20	35	
Married and completed child bearing	33	43	
Monthly income (CNY)			1.00
≤1000	113	115	
1000~3000	305	303	
3000~5000	41	40	
≥5000	49	49	
Degree pursuing			0.73
Graduate	385	389	
Postgraduate	123	118	
Double First-rate University			0.39
Yes	378	389	
No	130	118	
Types of students			0.40
Research oriented	280	266	
Professional oriented	228	241	
Kinds of research			0.28
Basic scientific research	194	169	
Clinical research	164	171	
Both	117	138	
Uncertain	33	29	
Daily research time (Hours)			0.64
1-5	215	233	
6-10	153	142	
11-15	126	116	
≥16	14	16	
Number of published academic paper			0.69
0	338	349	
1-2	119	111	
3-4	29	31	
≥5	22	16	
First author paper, SCI index			0.27
Yes=1	122	107	
No=0	386	400	
Number of participating projects			0.07
0	130	143	
1-3	357	355	
≥4	21	9	
Scientific learning style			0.34
Guiding by others	313	327	
Self-study	195	180	
Feelings of time stress (Range from 1-7, 1 is none, 2-3 mild, 4-5 moderate, 6-7 severe)			0.31
1 (None)	16	18	
2-3 (Mild)	45	52	

4-5 (Moderate)	159	132	
6-7 (Severe)	288	305	
Tutor is director of the department			0.59
Yes	289	280	
No	219	227	
Tutor is academician of the CAS and CAE			1.00
Yes	3	3	
No	505	504	
Tutor hold position in the national academic organizations			0.03
Yes	61	40	
No	447	467	
Tutor won a bid of NSFC within 5 years			0.44
Yes	356	344	
No	152	163	
GAD-7 (mean)	8.75±5.41	8.85±5.39	0.77
PHQ-9 (mean)	8.65±6.17	8.90±5.94	0.29
Severe mental distress			0.43
Yes	76	85	
No	432	422	
Abbreviation: SCI, Science Citation Index; CNY, Chinese Yuan; CAS, Chinese Academy of Sciences; CAE, Chinese Academy of Engineering; NSFC, National Natural Science Foundation of China; GAD-7, Generalized Anxiety Disorder Scale-7; PHQ-9, Patient Health Questionnaire-9. Notes: Severe mental distress was defined as the sum of GAD-7 and PHQ-9 scores ≥30.			

Table 2. The significant characteristics included in the model and corresponding estimates.		
Characteristics	Estimate ¹	Assigned score
Intercept	-5.06	
Year of study	0.20	
First year		1
Second year		2
Third year		3
Fourth year		4
In deferment period		5
Type of student	0.44	
Research oriented		1
Professional oriented		2
Daily research time (Hours)	0.51	
1-5		1
6-10		2
11-15		3
≥16		4
Monthly income (CNY)	-0.28	
≤1000		1
1000~3000		2
3000~5000		3
≥5000		4
Scientific learning style	0.30	
Guiding by others		1
Self-study		2
Feelings of time stress (Range from 1-7, 1 is none, 2-3 mild, 4-5 moderate, 6-7 Severe)	0.39	
1 (None)		1
2-3 (Mild)		2
4-5 (Moderate)		3
6-7 (Severe)		4
Abbreviation: CNY, Chinese Yuan. Notes: Seven predictors were identified by the LASSO method to be associated with the prevalence of severe mental distress significantly. Estimates were calculated using the logistic regression model. Severe mental distress was taken as Y in the model. Six variables were included in the model according to the coefficient values (value≤0.01). The formula: $P(Y=1) = \frac{e^x}{(1+e^x)}$, $x = -5.06 + 0.20 * \text{Year of study} + 0.44 * \text{Type of student} + 0.51 * \text{Daily research time} - 0.28 * \text{Monthly income} + 0.30 * \text{Scientific learning style} + 0.39 * \text{Feelings of time stress}$. $P(Y=1)$ indicates the predicted probability of severe mental distress.		

Table 3. The effective performances of the model in validation group.

Discrimination ability	AUROC	Slope ¹	95% CI	CCR	Sensitivity	Specificity	False POS	False NEG
Training group	0.70	0.06	0.04~0.08	82.30%	11.80%	94.70%	71.90%	14.10%
Validation group	0.66	0.04	0.02~0.06	81.30%	4.70%	96.70%	77.80%	16.60%
Calibration ability	Slope ²	95% CI	X-intercept	95% CI	Y-intercept	95% CI	P ³	
Training group	0.90	0.65~1.15	-0.02	-0.09~0.03	0.02	-0.03~0.06	0.82	
Validation group	0.80	0.51~1.09	-0.01	-0.11~0.04	0.01	-0.05~0.06	0.97	

Abbreviations: AUROC, the area under the receiver operating characteristic curve; CI, confident interval; CCR, correct classification rate; POS indicates positive; NEG indicates negative.
1 indicates discrimination slope;
2 indicates calibration slope;
3 indicates Hosmer and Lemeshow Goodness-of-Fit Test.

Table 4. Risk group classifications of the developed model.

Groups	Number of students	Predicted value P(Y=1)	Actual value P(Y=1)	p ¹
Low risk (0.00%~9.99%)	345	6.81%	6.96%	0.001
Moderate risk (10.00%~19.99%)	461	14.54%	16.49%	
High risk (20.00%~)	209	27.45%	29.19%	

Notes: P(Y=1) indicates the rate of severe mental distress.
1 indicates actual rate of severe mental distress among the three risk groups.

Figures

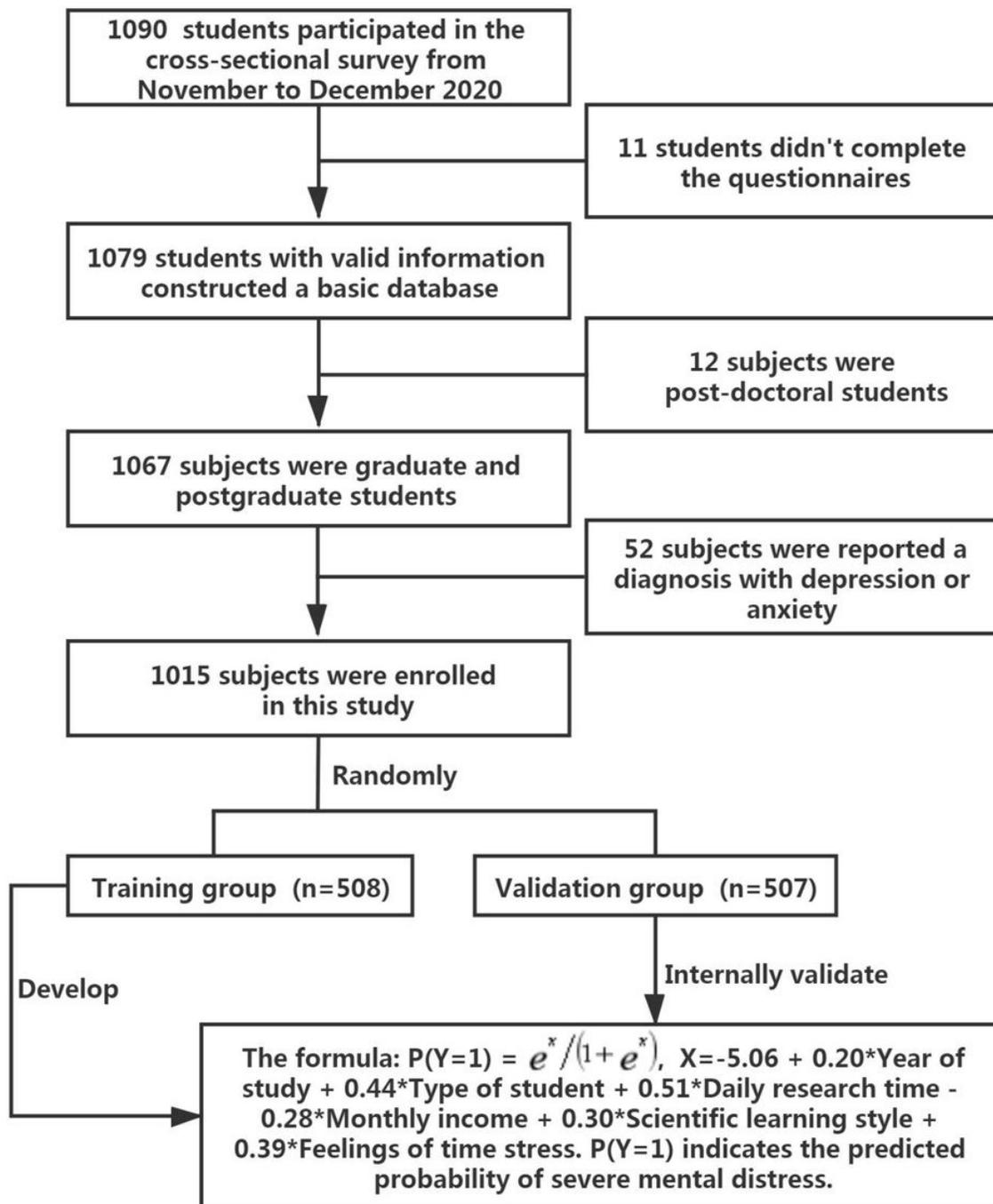


Figure 1

Study profile.

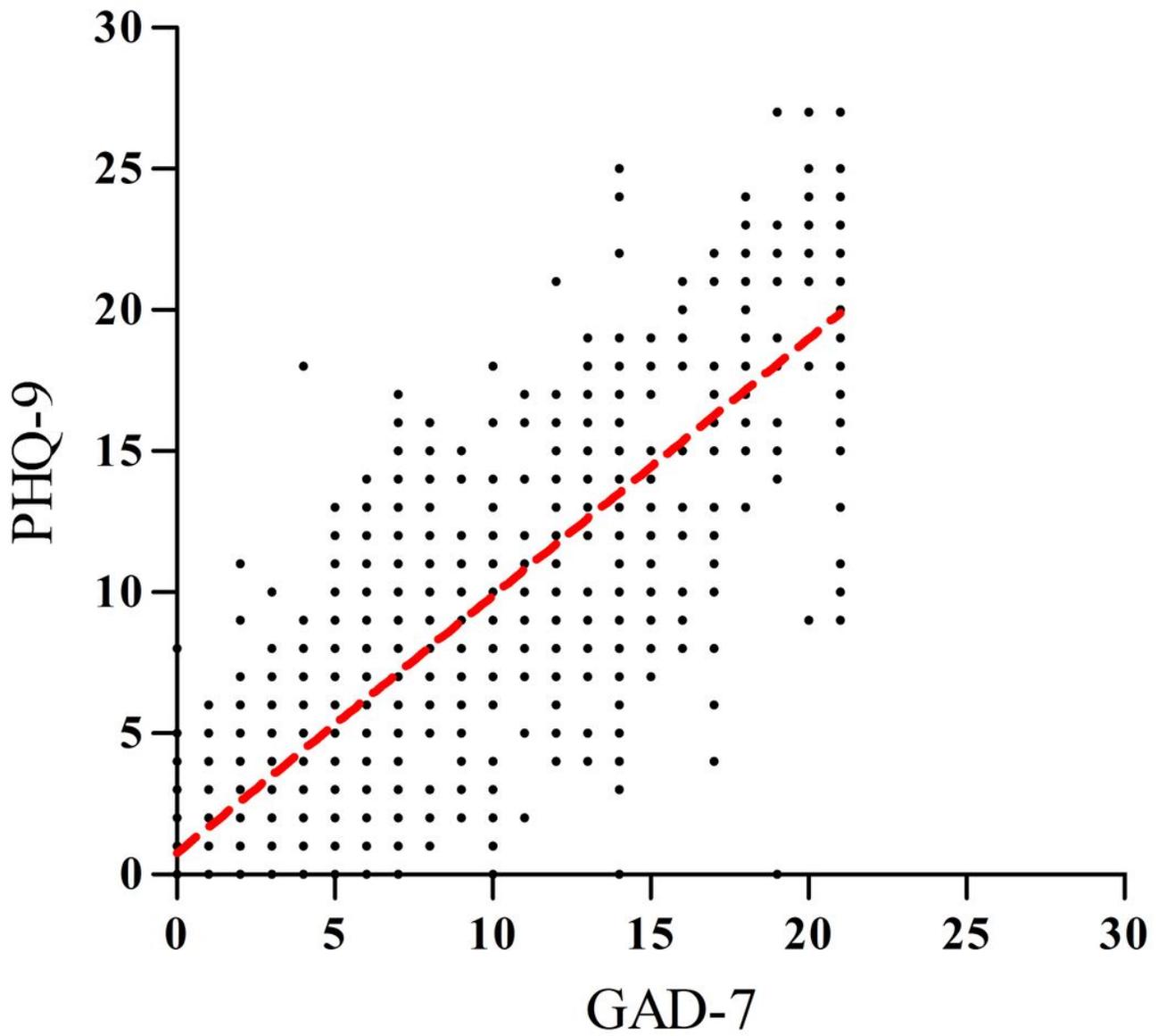


Figure 2

There is a positive correlation between GAD-7 and PHQ-9, and the curve can be fitted with $y=0.91x+0.75$ (GAD-7 is x, PHQ-9 is y).

ROC Curve for Model

Area Under the Curve = 0.6979

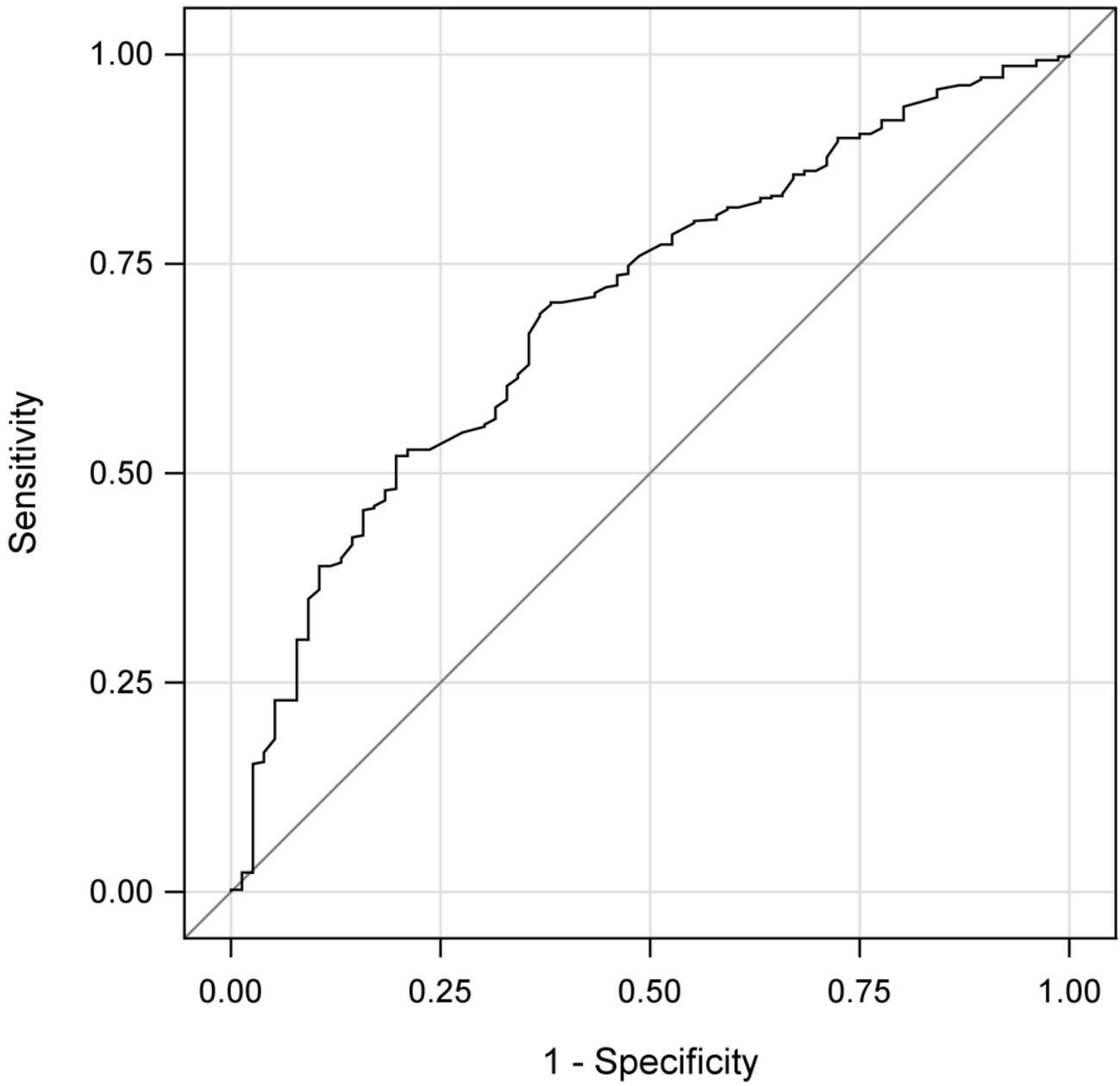


Figure 3

ROC curve for the developed model in the training group.

ROC Curve for Model

Area Under the Curve = 0.6620

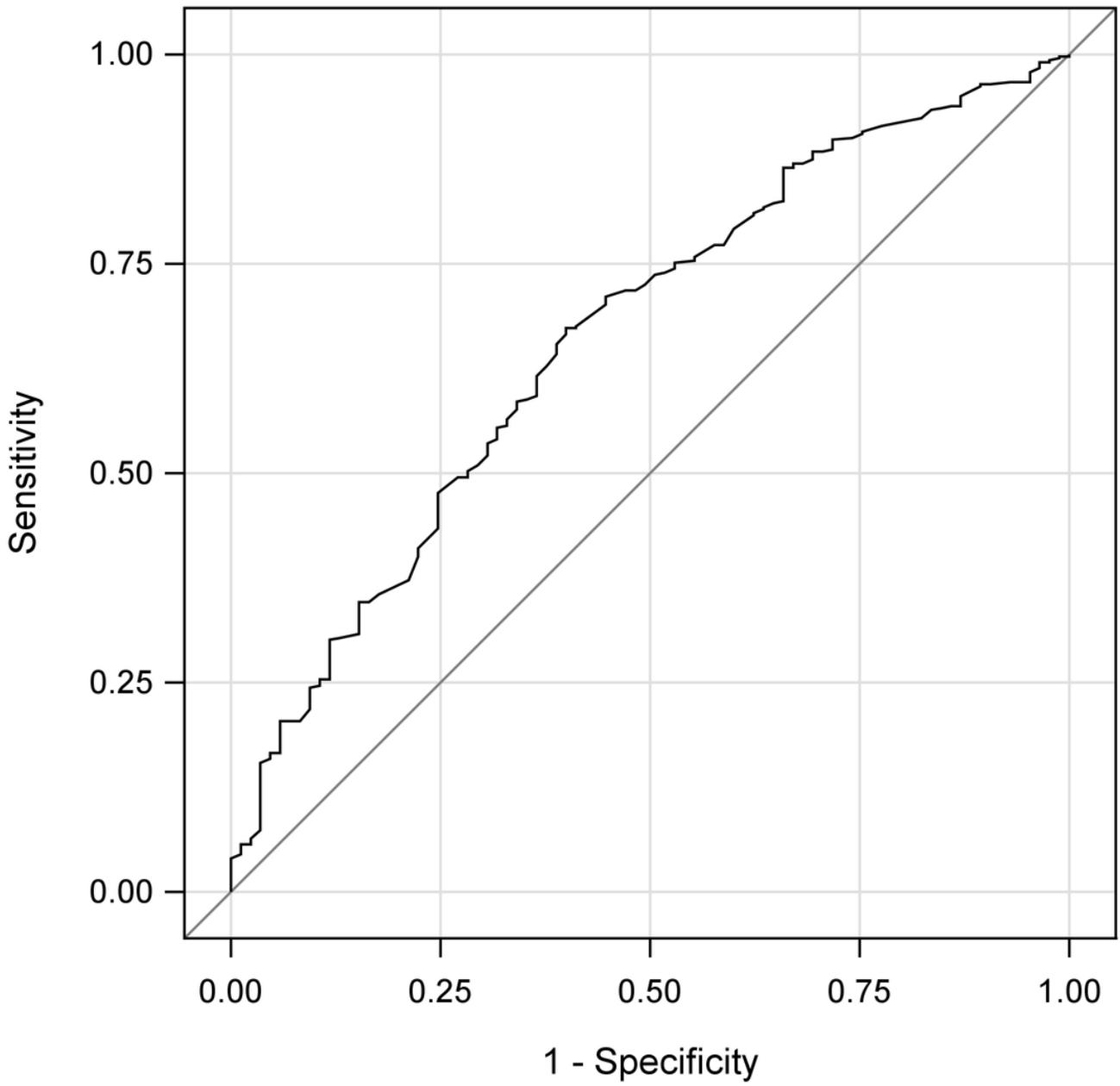


Figure 4

ROC curve for the developed model in the validation group.

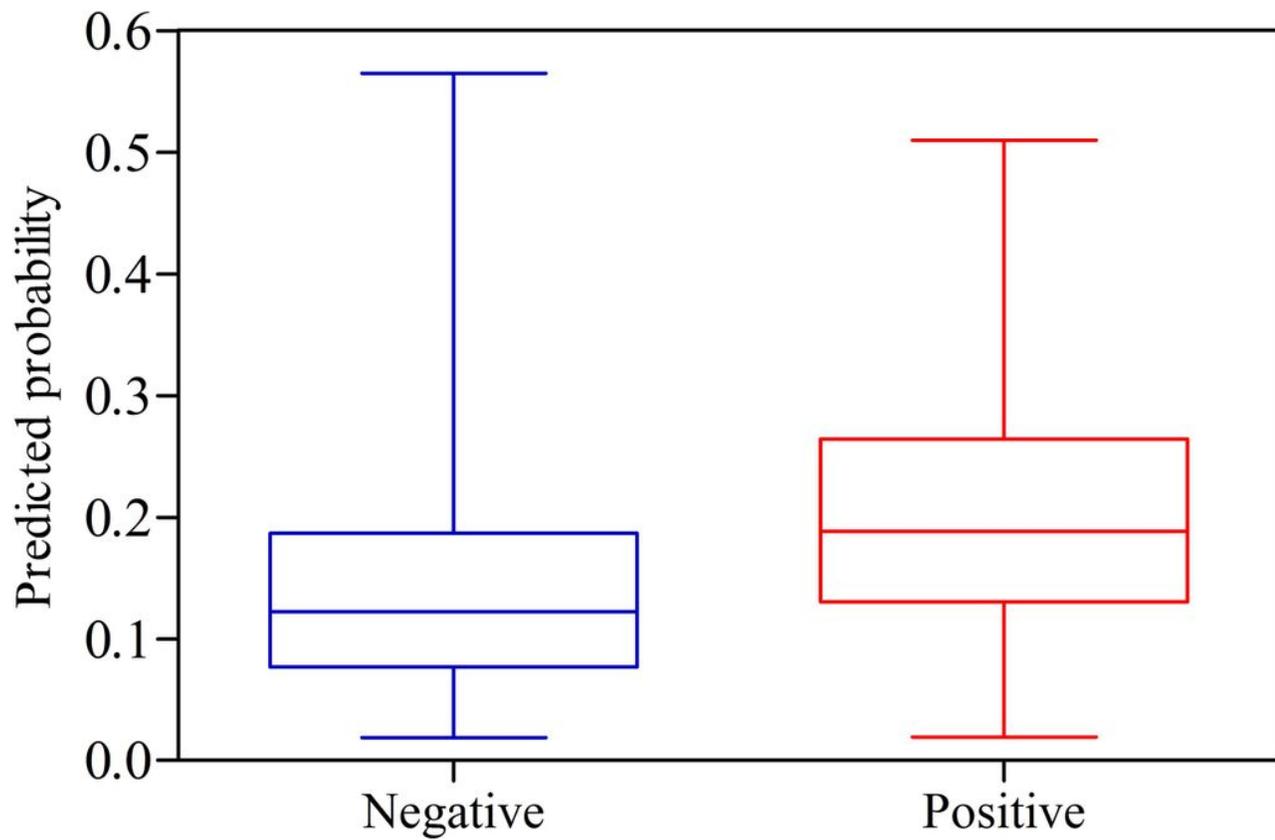


Figure 5

Discrimination slope of the developed model in the training group (Slope=0.06, Negative \approx 0.14 vs Positive \approx 0.20, $P \approx$ 0.001).

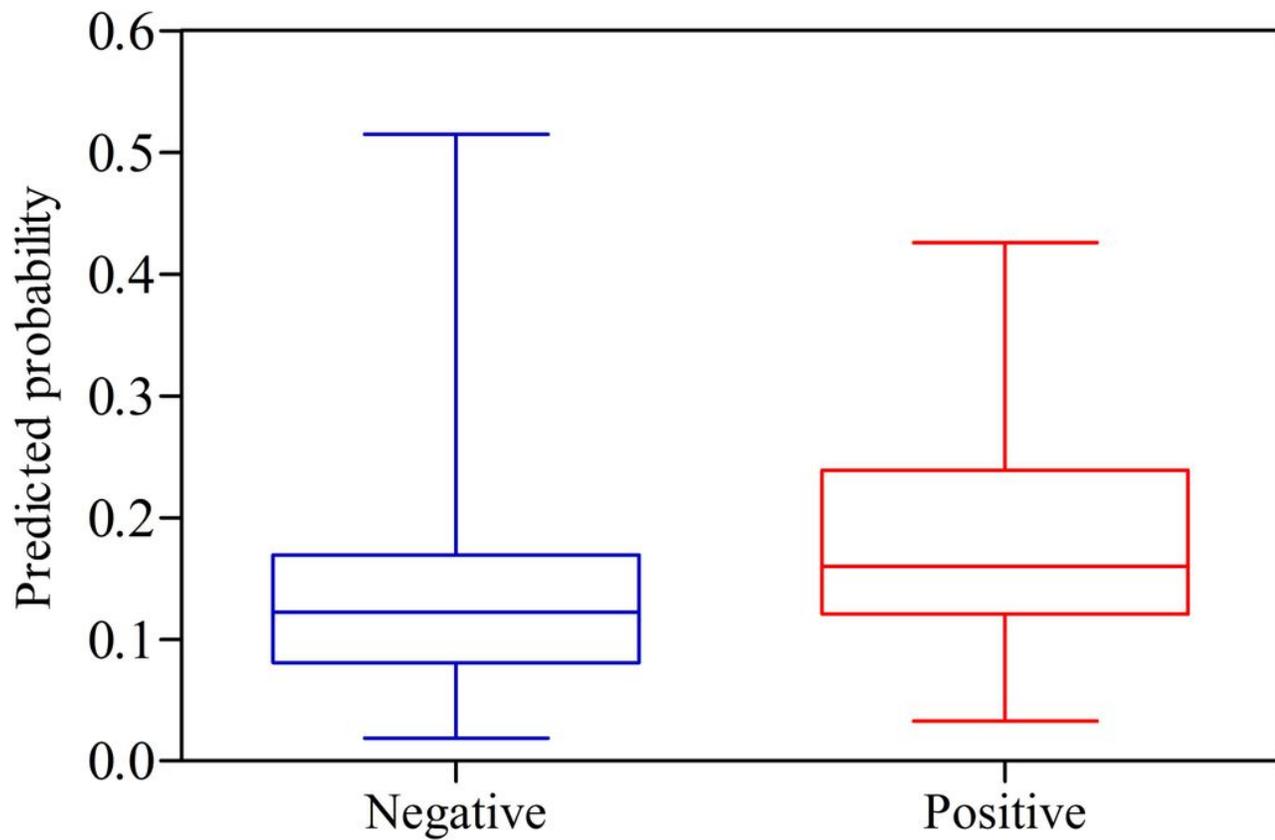


Figure 6

Discrimination slope of the developed model in the validation group (Slope=0.04, Negative \bar{x} 0.13 vs Positive \bar{x} 0.17, $P < 0.001$).

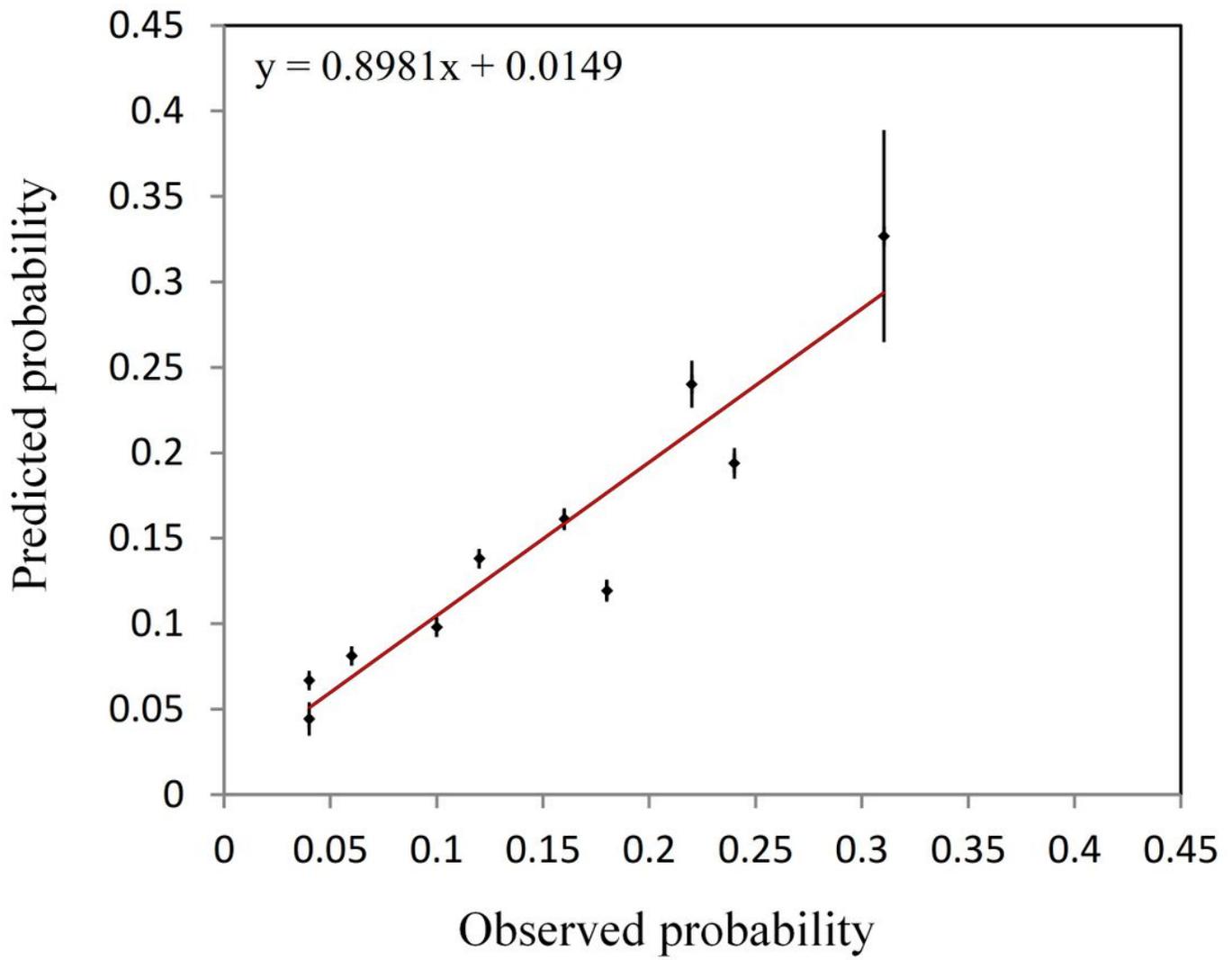


Figure 7

Calibration slope of the developed model in the training group (Slope=0.90, $P < 0.001$).

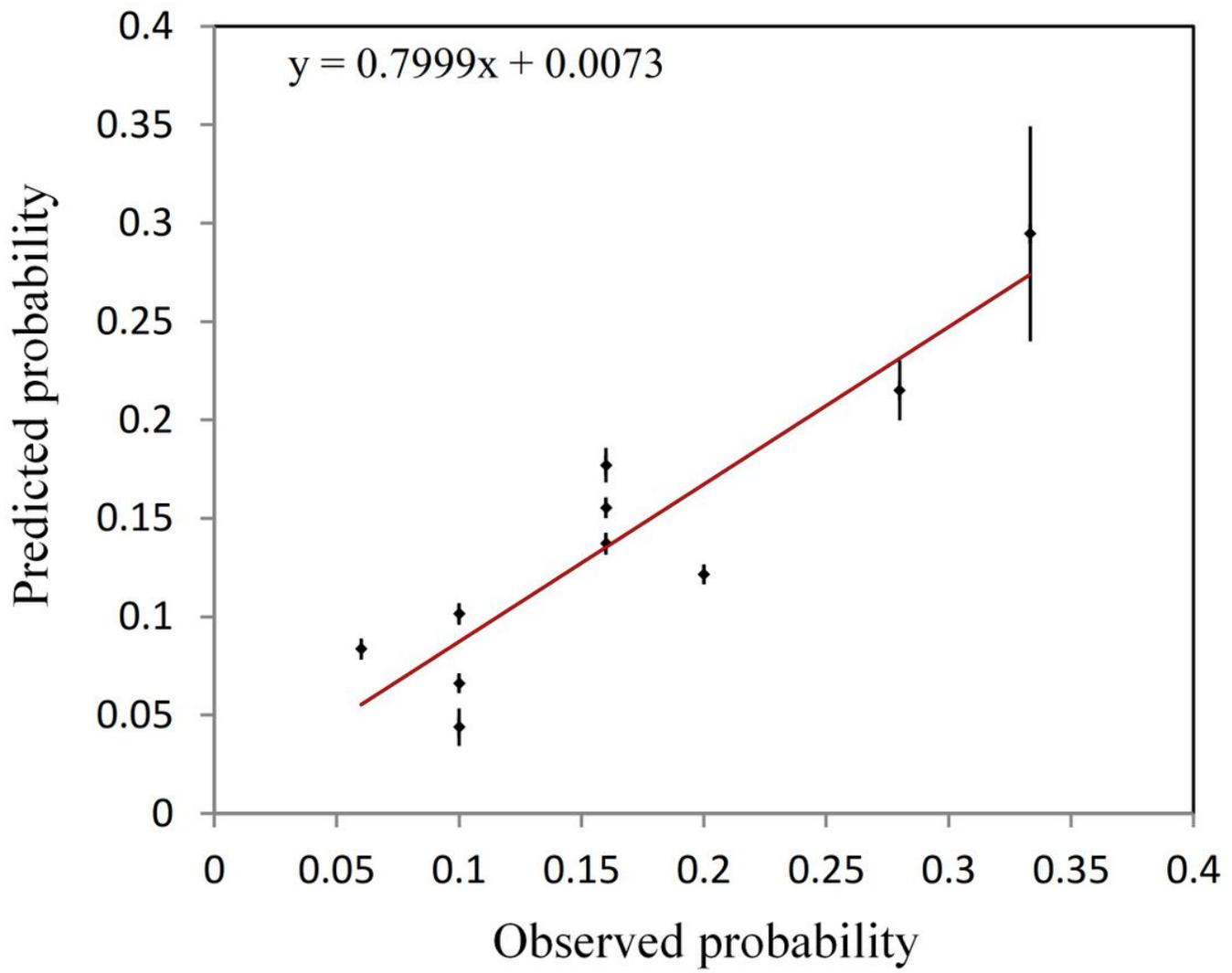


Figure 8

Calibration slope of the developed model in the validation group (Slope=0.80, $P < 0.001$).