

Machine Learning Based Early Diagnosis of ADHD with SHAP Value Interpretation: A Retrospective Observational Study

Xinyu Zhang^{1,*}, Xue Xiao^{1,*}, Yufan Luo¹, Wei Xiao¹, Yingsi Cao¹, Yuanjin Chang¹, Dongqin Wu¹, Hua Xu¹, Jinlin Zhao¹, Xianhui Deng², Yuanying Jiang³, Ruijin Xie^{1,4}, Yueying Liu¹ 

¹Department of Pediatrics, Affiliated Hospital of Jiangnan University, Wuxi, People's Republic of China; ²Department of Neonatology, Jiangyin People's Hospital of Nantong University, Wuxi, People's Republic of China; ³Linping Campus, The Second Affiliated Hospital of Zhejiang University School of Medicine, Hangzhou, People's Republic of China; ⁴Yangzhou Polytechnic College, Yangzhou, People's Republic of China

*These authors contributed equally to this work

Correspondence: Ruijin Xie; Yueying Liu, Department of Pediatrics, Affiliated Hospital of Jiangnan University, Wuxi, People's Republic of China, Email xieruijin@gmail.com; shadow7671@163.com

Background: Attention-Deficit/Hyperactivity Disorder (ADHD) is a common neurodevelopmental disorder in children, characterized by inattention, hyperactivity, and impulsivity. Current diagnostic methods for ADHD rely primarily on behavioral assessments, which can be challenging due to symptom overlap with other psychiatric disorders and significant inter-individual variability. Developing potential early diagnostic methods for ADHD is imperative to mitigate the risk of misdiagnosis and enhance the evaluation of treatment efficacy.

Methods: The study was conducted at the Department of Pediatrics, Affiliated Hospital of Jiangnan University, from November 2022 to January 2024. Clinical data, including complete blood count, liver and kidney function tests, blood glucose levels, serum electrolyte tests, and serum 25-dihydroxyvitamin D3 levels, were collected. Feature selection and model construction were performed using various machine learning algorithms.

Results: Our results indicated that the Gradient Boosting Machine algorithm is the optimal model.

Conclusion: Our machine learning analyses suggest that the Gradient Boosting Machine (GBM) model may be the optimal choice, highlighting blood beta-2 microglobulin levels, red blood cell distribution width, 25-dihydroxyvitamin D3, and the percentage of eosinophils as key predictors of ADHD risk, thereby aiding early diagnosis. Further large-scale studies are warranted to validate these findings and explore the underlying mechanisms.

Keywords: ADHD, diagnosis, biomarkers, machine learning, SHAP methods

Introduction

Attention-Deficit/Hyperactivity Disorder (ADHD) is a prevalent neurodevelopmental disorder in children, typically manifesting before the age of 12.¹ It is characterized by symptoms of inattention, hyperactivity, and impulsivity, with a higher prevalence in males than in females.² ADHD can lead to academic challenges, impaired social relationships, and diminished self-esteem. Furthermore, recent data from the National Health Interview Survey (NHIS) for the years 2020–2022 indicate that the prevalence of ever-diagnosed ADHD in children aged 5–17 years in the United States is 11.3%.³ In China, the burden of ADHD is expected to increase due to sociodemographic transitions and growing awareness of diagnostic criteria. Studies report that the prevalence of ADHD among Chinese children ranges from 4.96% to 9.8%, with considerable regional variability. For example, a study conducted in Deyang, Sichuan Province, found a prevalence of 9.8% among primary school students, while another study in rural areas of China reported a prevalence of 7.5%. These findings underscore the importance of early diagnosis and intervention to help children manage their symptoms and improve their overall quality of life.

Currently, ADHD diagnosis primarily relies on behavioral assessments and the criteria outlined in the Diagnostic and Statistical Manual of Mental Disorders, Fifth Edition (DSM-5).⁴ Early diagnosis in clinical settings remains challenging due to various factors. ADHD symptoms often overlap with those of other psychiatric disorders, complicating differential diagnosis.⁵ Moreover, ADHD symptoms exhibit significant inter-individual variability and can manifest differently across developmental stages.⁵ Additionally, comorbid conditions, such as learning disabilities and mood disorders, can further complicate the diagnostic process.⁶ Furthermore, early diagnosis and intervention can significantly improve a child's academic performance, social relationships, and overall well-being.⁷ By addressing ADHD symptoms early, individuals can develop coping strategies and skills to manage their challenges effectively.⁸ Consequently, developing potential early diagnostic methods for ADHD, particularly using serum biomarkers, is imperative to mitigate the risk of misdiagnosis and enhance the evaluation of treatment efficacy.

Recently, early diagnosis of neurodevelopmental disorders, including ADHD, using machine learning-based predictions hold immense potential for transforming how we understand and support individuals with these conditions. Machine learning algorithms are proving to be powerful tools for analyzing complex datasets, identifying patterns indicative of neurodevelopmental disorders, and helping physicians identify shared and distinct features that lead to a better understanding of these conditions' underlying mechanisms.⁹ However, there are few studies focusing on the early diagnosis of ADHD using machine learning-based approaches. Therefore, this study aims to establish a machine learning-based risk probability predictive model for ADHD to aid in the early diagnosis of ADHD by integrating observational cohort studies and advanced machine learning algorithms. For this purpose, we collected data including complete blood counts, liver function assessments, kidney function evaluations, blood glucose levels, serum electrolyte analyses, and vitamin D level. These indicators were selected not only because they are widely accessible across all healthcare facility levels in China but also due to their high cost-effectiveness, making them particularly suitable for large-scale screening of ADHD in China.¹⁰ This study will not only enhance our understanding of the biological basis of ADHD but may also improve treatment strategies for timely intervention and management.

Methods

Study Design

As shown in [Figure 1](#), this pilot prospective observational cohort study was conducted at the Department of Pediatrics, Affiliated Hospital of Jiangnan University, from 21 November 2022 to 12 January 2024, in accordance with our previous study.¹¹ Children aged 1 to 18 years diagnosed with ADHD were considered for the ADHD group, while healthy children admitted to the hospital for routine check-ups were included in the control group. The basic clinical characteristics of both groups are detailed in [Table 1](#). We excluded children who with digestive system disorders or other mental disorders, and those with missing clinical or laboratory data. To facilitate the identification of potential serum biomarkers for ADHD in clinical settings, we focused on differences in complete blood count, liver function tests, kidney function tests, blood glucose levels, serum electrolyte tests, and serum 25-dihydroxyvitamin D3 levels. These data were collected within the first 24 hours of admission to the department. Initially, a total of 50 features were investigated ([Supplemental Table 1](#)).

Ethical Statement

This preliminary observational pilot study was conducted in accordance with the principles of the Declaration of Helsinki and received ethical approval from the Institutional Review Board of the Affiliated Hospital of Jiangnan University (ID: LS202112R2). Participants, as well as their parents or legal guardians, were comprehensively informed about the study's purpose, procedures, and potential risks and benefits. Written informed consent was obtained from all participants and their parents or legal guardians prior to their involvement. Participation was entirely voluntary, and participants were advised of their right to withdraw at any time without repercussions. Participation was entirely voluntary, and participants were advised of their right to withdraw at any time without repercussions. All collected data were anonymized and securely stored to ensure participant confidentiality, preventing the identification of individual participants during data collection.

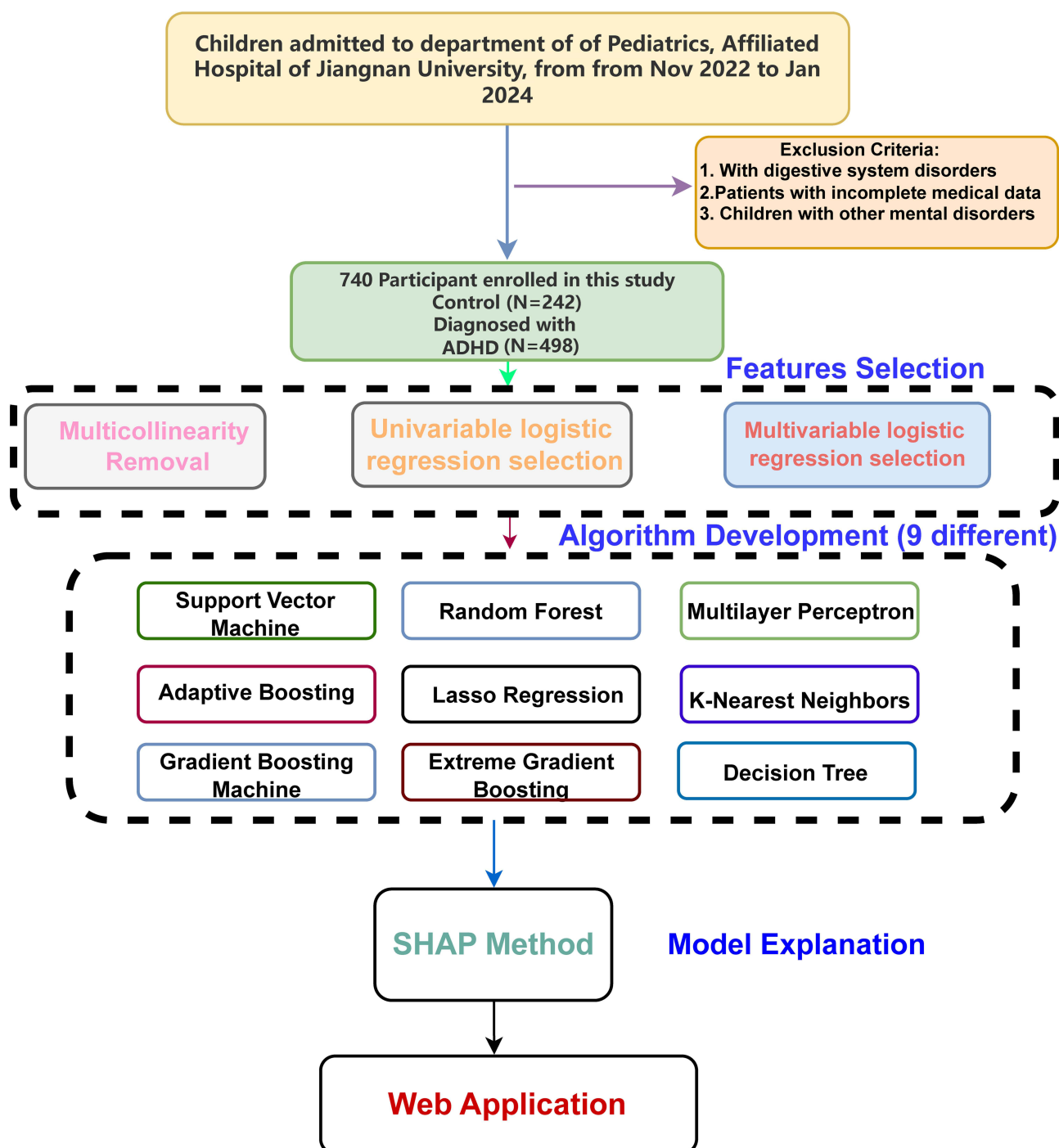


Figure 1 Flow chart of the study design.

Clinical Diagnostic Criteria for ADHD

The clinical diagnostic criteria for Attention-Deficit/Hyperactivity Disorder (ADHD) are based on the DSM-5 (Diagnostic and Statistical Manual of Mental Disorders, Fifth Edition).¹² In brief, six or more symptoms from either or both categories of inattention and hyperactivity-impulsivity for children up to age 16, or five or more for adolescents aged 17 and older and adults, must be present for at least 6 months.¹³ These symptoms must be inappropriate for the individual's developmental level and interfere with daily functioning.¹³ Importantly, several symptoms must have been

Table 1 The Basic Clinical Characteristics of Participants Involved in This Study

Characteristic	Overall, N = 740	Control, N = 242	ADHD, N = 498	P-value [#]
Age, Mean (SD)	8.95 (2.46)	8.74 (2.91)	9.05 (2.20)	0.92
Sex, n (%)				0.74
Male	559 (76)	181 (75)	378 (76)	
Female	181 (24)	61 (25)	120 (24)	
Education, n (%)				<0.001
Kindergarten	31 (4.2)	27 (11)	4 (0.8)	
Primary school	556 (75)	152 (63)	404 (81)	
Junior high school	153 (21)	63 (26)	90 (18)	
BMI, Mean (SD)	16.84 (1.70)	16.87 (1.63)	16.82 (1.74)	0.73
Birth type, n (%)				0.19
Vaginal Delivery	670 (91)	224 (93)	446 (90)	
Cesarean Section	70 (9.5)	18 (7.4)	52 (10)	

Notes: [#]Wilcoxon rank sum test; Pearson's Chi-squared test.

Abbreviation: BMI, Body Mass Index.

present before the age of 12 years. The detailed information of clinical diagnostic criteria for ADHD, including specific symptom lists and exclusion criteria, are shown in [Supplemental Method 1](#).

Feature Selection

As shown in [Figure 1](#), to mitigate multicollinearity bias, we performed Spearman correlation analysis and excluded variables with significant correlations ($R > 0.9$), in accordance with the methodology outlined in the aforementioned study initially.¹⁴ Next, we performed feature selection using both univariable and multivariable logistic regression analyses, based on methods from our previous studies.¹⁵ Additionally, we constructed a traditional risk probability predictive model for ADHD using a nomogram, and verified the discrimination and efficacy of our feature selection by comparing the receiver operating characteristic (ROC) curve, decision curve analysis (DCA), and calibration curve, as described in a previous study.¹⁴

Model Construction and Evaluation

As shown in [Supplemental Table 2](#), a total of 15 features were selected for the development of the prediction models. We employed five different machine learning (ML) models to predict the risk of ADHD: 1. Adaptive Boosting (AdaBoost), 2. Lasso Regression (Lasso), 3. Random Forest (RF), 4. Gradient Boosting Machine (GBM), 5. Extreme Gradient Boosting (XGBoost), 6. Support Vector Machine (SVM), 7. K-Nearest Neighbors (KNN), 8. Multilayer Perceptron (MLP), and 9. Decision Tree (DT). The use of multiple models enabled us to compare their performance and identify the most effective approach for predicting the risk of ADHD. To evaluate the reliability of these models, we used three primary metrics: the area under the Receiver Operating Characteristic (ROC) curve (AUC), the Precision-Recall Curve (PRC), and Decision Curve Analysis (DCA).¹⁴

Model Interpretation

The interpretability of machine learning has always been challenging.¹⁷ To further explain how each feature variable affects and contributes to the final model, we employed the SHapley Additive eXplanation (SHAP) method to interpret the highest-performing black-box model based on previous studies.¹⁸ In this study, we evaluated the importance of each feature by computing the mean absolute SHAP value. We also plotted the SHAP values for each feature across samples to better understand the overall patterns and the impact range of features on the risk of ADHD. We also utilized the SHAP dependency plot to evaluate the effects of each feature. Additionally, we provided two examples of SHAP predictions for demonstration purposes. Furthermore, to facilitate the utility of the model in clinical settings, the final prediction model was implemented into a web application developed using the Streamlit Python framework and available at <https://adhdrisk.streamlit.app/>.¹⁹

Table 2 Detailed Information of GWAS Data Used in the Study

GWAS ID	Years	Trait	Sample Size	Population	Database/PMID*
ieu-a-1183	2017	ADHD	55,374	European	IEU Database
ebi-a-GCST90025948	2021	Serum phosphate levels	400,159	European	IEU Database PMID: 34226706
ebi-a-GCST90025980	2021	Aspartate aminotransferase levels	436,275	European	IEU Database PMID: 34226706
prot-c-3485_28_2	2019	Blood beta-2-microglobulin	3,080	European	IEU Database PMID:28240269
ebi-a-GCST004606	2017	Blood eosinophil counts	172,275	European	IEU Database PMID:27863252
ebi-a-GCST90018973	2021	Total bilirubin levels	34,829	European	IEU Database PMID:34594039
ebi-a-GCST90000615	2020	Vitamin D level	417,580	European	IEU Database PMID: 32242144
ukb-d-30070_irnt	2018	Red blood cell distribution width	350,473	European	IEU Database
ebi-a-GCST90025992	2021	Serum albumin levels	400,938	European	IEU Database PMID: 34226706

Notes: *IEU OPEN GWAS Database (<https://gwas.mrcieu.ac.uk/>).¹⁶

Mendelian Randomization Analysis

Mendelian Randomization (MR) analysis is a powerful technique that leverages genome-wide association studies (GWAS) data and genetic variants to explore potential relationships between modifiable exposures and various outcomes or diseases.²⁰ GWAS identifies genetic variants associated with specific traits, diseases, or outcomes, while genetic variants, especially single-nucleotide polymorphisms (SNPs), serve as instrumental variables (IVs) in MR studies.²¹ MR analysis has been widely used to assess the potential effects of various exposures on disease risk.²¹ In this study, to explore potential association between ADHD and potential serum biomarkers, we conducted bi-directional Mendelian Randomization analysis based on previous studies.^{22–24} [Supplemental Method 2](#) provides detailed information about the criteria for Mendelian randomization analysis, and [Table 2](#) shows the detailed information of GWAS data used in this study. All GWAS data involved in this study were obtained from the IEU OpenGWAS project (<https://gwas.mrcieu.ac.uk/>).²⁵ We then utilized the SNPnexus database (<https://www.snp-nexus.org/v4/>), a web-based tool that provides an aggregate set of functional annotations for SNPs, to map SNPs to their closest genes based on IVs of potential serum biomarkers (detailed SNPs and associated genes are shown in [Supplemental Table 3](#)).²⁶ And the analysis of Gene Ontology (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway enrichment using the R package “clusterProfiler”,^{27,28} were conducted based on our preceding work to explore the biological pathways involving the pathophysiology of ADHD.^{22,24,29}

Statistical Analysis

In this study, data are presented as mean ± standard deviation (SD). The normality of the data was assessed using the Shapiro–Wilk test. For data following a normal distribution, statistical comparisons were conducted using unpaired t-tests, one-way ANOVA, or two-way ANOVA, followed by Tukey’s post-hoc tests for multiple comparisons. Non-normally distributed data were analyzed using the Mann–Whitney *U*-test for two groups or the Kruskal–Wallis test for multiple groups. All statistical analyses were performed using Python version 3.11.9 (<https://www.python.org>) and R Version 4.4.1 (<https://cran.r-project.org/>). A P-value of less than 0.05 was considered to indicate statistical significance. The raw code for the study is available through our GitHub repository (<https://github.com/PediatricLab-YueyingLiu/ADHDM>).

Results

Data Processing Results and Features Selection

As shown in [Figure 1](#) and [Table 1](#), after selection, a total of 740 volunteer participants were involved in this study, and there were no significant differences in age, sex, BMI, and birth type between the two groups ($P > 0.05$). Although there was a significant difference in educational attainment between the two groups ($P < 0.05$), educational attainment itself does not directly affect the risk of developing ADHD. Recent research has indicated that socioeconomic status can influence both the risk of ADHD and educational attainment.³⁰ Children from lower socioeconomic backgrounds may have less access to resources and support, which can exacerbate the challenges associated with ADHD.³¹

Then, univariable and multivariable logistic regression analyses were performed for feature selection ([Figure 2A](#) and [B](#)). A total of 15 features, including serum 25-dihydroxyvitamin D, were selected from an initial pool of 50 features ([Supplementary File 1](#)). Notably, high odds ratios (ORs) for absolute eosinophil count were observed in both univariable and multivariable logistic regression analyses. This may be attributed to the fact that eosinophils are a type of white blood cell involved in the body's immune response, and research has suggested that immune system abnormalities, such as inflammation and allergic responses, could potentially influence brain function and behavior.³² Additionally, there is no significant correlation existed between the variables, as illustrated in [Supplementary File 1](#). The nomogram visually combines multiple factors, including serum 25-dihydroxyvitamin D, to deliver an overall risk assessment ([Figure 2C](#)). The receiver operating characteristic (ROC) curve validated the discrimination and efficacy of our feature selection compared to the single factors, with highest AUC scores in the nomogram model (0.87) ([Figure 2D](#)). The calibration curve further supports the discrimination and efficacy of our feature selection, as neither the bias-corrected nor the apparent line is far from the ideal line ([Figure 2E](#)). Finally, the decision curve analysis (DCA) supports the model's potential clinical utility by comparing the "Model" line to the "None" line and the "All" line ([Figure 2F](#)). In summary, after data processing and feature selection via univariable and multivariable logistic regression analyses, our study demonstrated that a total of 15 features could serve as a predictive model for the risk probability of ADHD to aid in the early diagnosis of ADHD.

Explainable Machine Learning Analyses Results

After univariable and multivariable logistic regression analyses, we applied several widely used machine learning algorithms, including Adaptive Boosting (AdaBoost), Lasso Regression (Lasso), Random Forest (RF), Gradient Boosting Machine (GBM), Extreme Gradient Boosting (XGBoost), Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Multilayer Perceptron (MLP), and Decision Tree (DT), to further develop a machine learning-based predictive model for the probability of ADHD risk. As shown in [Figure 3A–C](#), we evaluated the performance of these predictive models using the receiver operating characteristic (ROC) curve, the precision-recall curve (PRC), and the decision curve analysis (DCA). Among these 9 different machine learning algorithms, the GBM model demonstrated the highest performance in the ROC (0.91), PRC (0.95), and DCA evaluations. Therefore, we chose the GBM model as the optimal model for this study. Additionally, we determined the optimal number of features among the machine learning algorithms by comparing the AUC scores of the ROC curves for different machine learning models with varying feature numbers ([Figure 3D](#)). The results indicated that 8 features constituted the optimal number for GBM model.

To facilitate interpretability, we applied the Shapley Additive Explanations (SHAP) method to the GBM-based predictive model used in this study. This method illustrates how the features affect the model's output (ADHD risk), as indicated in a previous study.¹⁴ The summary bar plot ([Figure 3E](#)) shows the eight evaluated risk factors based on their SHAP values, with red and blue dots in each feature importance row representing high-risk and low-risk values, respectively. The summary dot plot ([Figure 3F](#)) also displays the important features and their ranking, with both plots highlighting the key role of top 4 features including levels of blood beta 2 microglobulin, red blood distribution width, 25-dihydroxyvitamin D3, and percentage of eosinophil in predicting risk probability of ADHD. The SHAP dependence plot ([Figure 4A](#)) was used to understand how individual features affect the GBM model's output, with SHAP values higher than zero indicating a higher risk of ADHD. Two typical examples were provided, one for a low-risk (healthy) individual ([Figure 4B](#)) and another for a high-risk (ADHD) individual ([Figure 4C](#)). These examples suggest that vitamin D levels may be the key feature in lowering the risk probability of ADHD, as they tend to decrease the risk in both low-risk and high-risk scenarios. The SHAP Force Plot shows how a machine learning model arrives at a specific prediction for an individual instance ([Figure 4D](#)). The final prediction model was implemented into a web application to facilitate its utility in

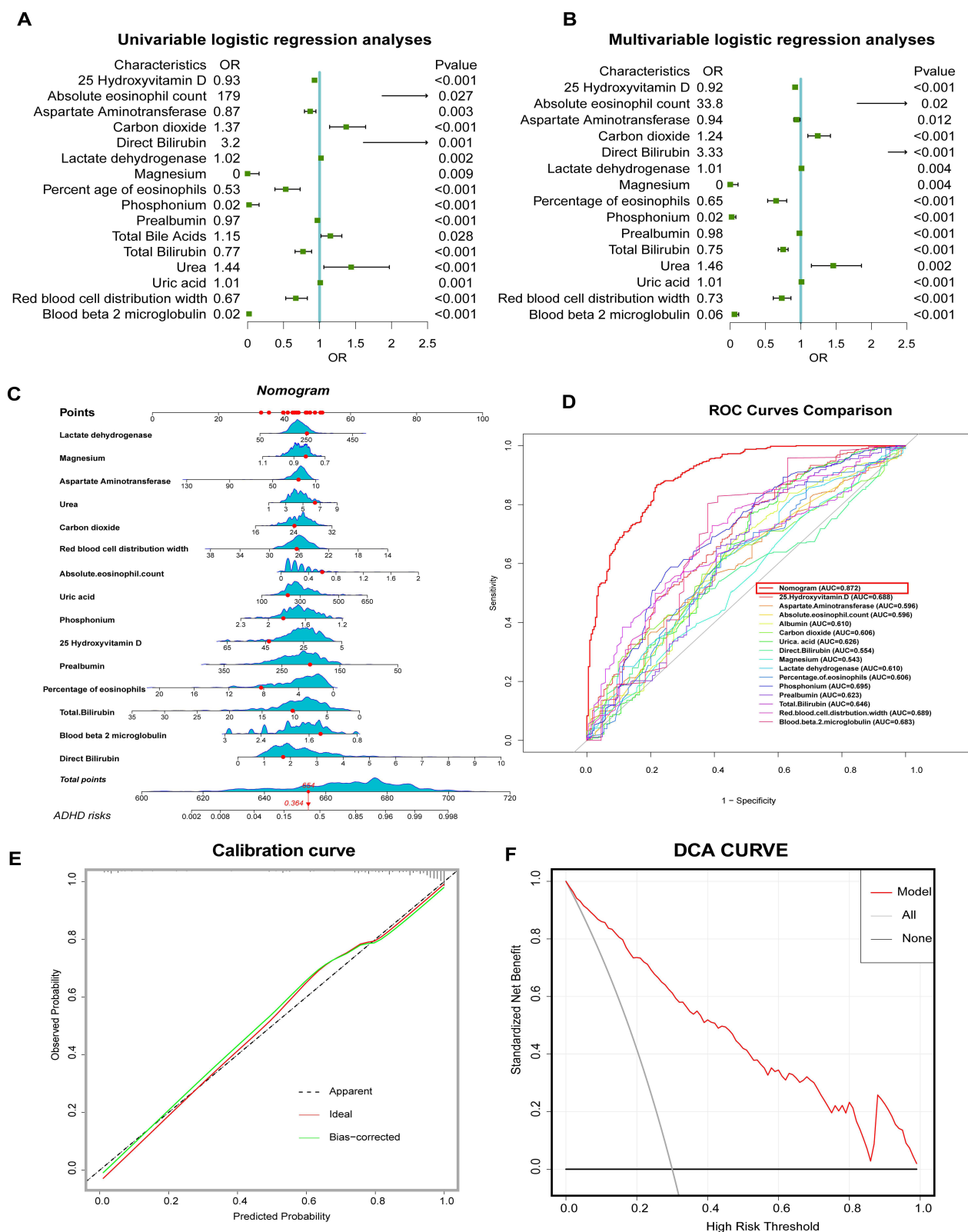


Figure 2 Univariable and multivariable logistic regression analyses. **(A and B)** Univariable and multivariable logistic regression analyses. **(C)** Nomogram is conducted to predict the risks probability of ADHD. **(D)** The ROC (Receiver Operating Characteristic) curve. **(E)** The Calibration curve. The closer the Bias-corrected or Apparent line is to the Ideal line, the better the model's calibration. **(F)** DCA (Decision Curve Analysis) curve. The "None" line represents the net benefit of not using any prediction model and the line labeled "All" represents the net benefit of treating all patients as high risk, regardless of their actual risk level. By comparing the "Model" line to the "None" line or "All" line, the DCA curve help evaluate the clinical usefulness of a predictive model.

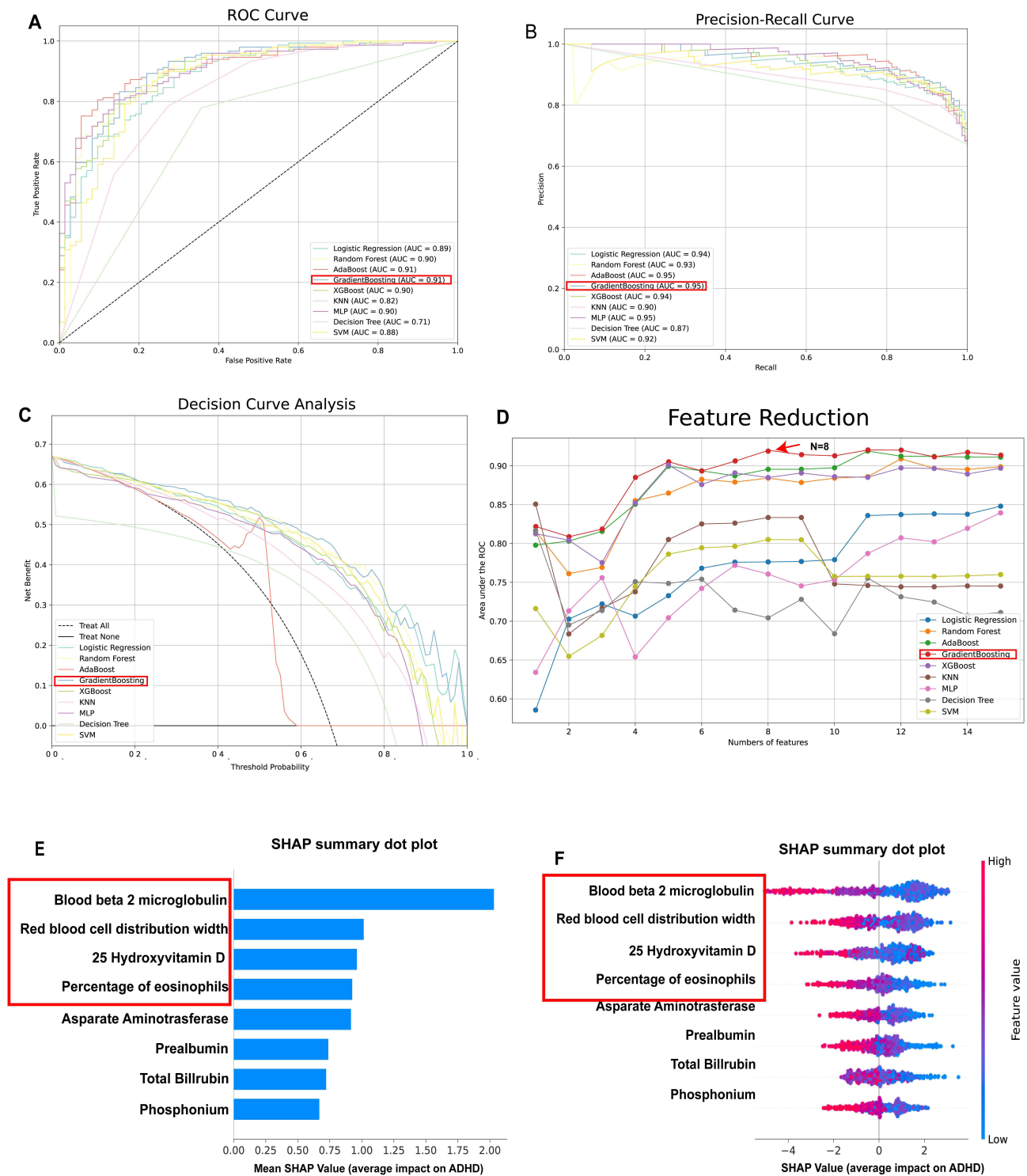


Figure 3 Comprehensive evaluation of machine learning models. **(A–C)** ROC (Receiver Operating Characteristic) curve, Precision-Recall (PRC) curve and DCA (Decision Curve Analysis) curve of different machine learning models. **(D)** AUC score in ROC curve of different machine learning models with varied feature numbers. **(E)** Summary bar plot by the SHAP method, showing the average impact of each feature on the risk of ADHD, with features sorted by importance. **(F)** Summary dot plot by the SHAP method, providing a more detailed view of each feature's impact. Each dot represents a single participant, with color indicating the feature value (red for high, blue for low) and horizontal position indicating the impact on the risk of ADHD. **(F)** SHAP summary dot plot illustrating the overall distribution of each feature's influence, where red generally indicates high values and blue indicates low values for the corresponding feature.

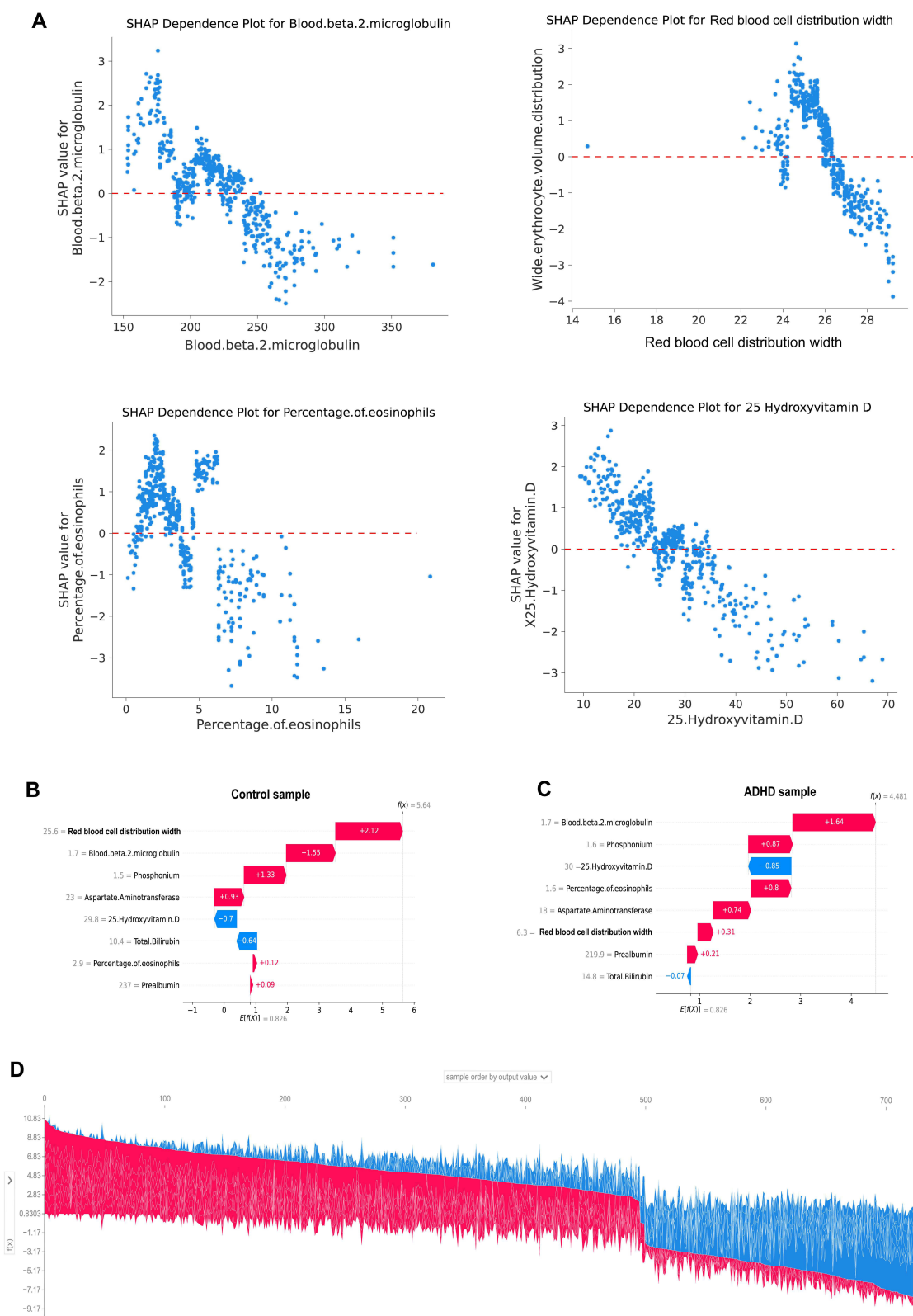


Figure 4 (A) The SHAP dependence plot of TOP 4 features selected by the GBM model, deemed most important for predicting ADHD risk. These plots show the relationship between individual feature values and their impact on the model's output (ADHD risk prediction). SHAP values (the y-axis) above zero push the model's decision towards predicting ADHD, while values below zero push the decision away from predicting ADHD. (B and C) Representative waterfall plots for a healthy child and an ADHD child. Red arrows indicate features that increase the risk of ADHD, while blue arrows indicate features that decrease the risk of ADHD. (D) Force plot for all participants. The x-axis represents individual participants, ordered by their predicted ADHD risk, and the y-axis represents the contributions of different features to the model's output.

Predictive Model for ADHD Risk Assessment

25.Hydroxy.vitamin.D

6.00 - +

Percentage.of.eosinophils

12.00 - +

Phosphonium

4.00 - +

Prealbumin

5.00 - +

Red.blood.cell.distribution.width

7.00 - +

Blood.beta.2.microglobulin

8.00 - +

Aspartate.Aminotransferase

11.00 - +

Total.Bilirubin

12.00 - +

Predict

Based on feature values, predicted possibility of ADHD is: 76.72%

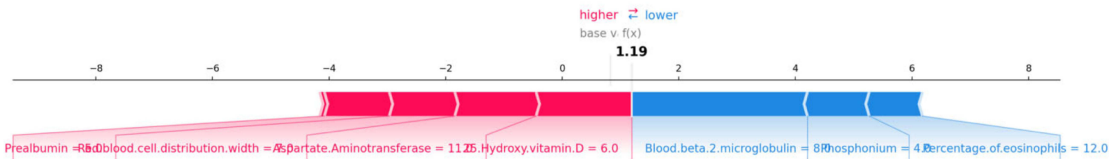


Figure 5 The website application to estimate ADHD risk. The user-friendly application of the final GBM model, which uses 8 features, is accessible for ADHD prediction.

clinical scenarios, as shown in Figure 5, and is available through <https://adhdrisk.streamlit.app/>. In summary, our machine learning analyses suggest that the GBM model may be the optimal choice, highlighting blood beta-2 microglobulin levels, red blood cell distribution width, 25-dihydroxyvitamin D3, and the percentage of eosinophils as key predictors of ADHD risk, thereby aiding early diagnosis.

Mendelian Randomization Analyses Results

To further explore and verify the association between the selected 8 features from the GBM predictive model and ADHD symptoms, bi-directional Mendelian randomization analyses were conducted based on our previous study.^{24,29} As shown in Figure 6A and B, bi-directional Mendelian randomization analyses indicated that only serum 25-dihydroxyvitamin D3 was associated with ADHD symptoms in both directions ($P < 0.05$), and ADHD symptoms may influence total bilirubin levels (Figure 6B). Furthermore, we performed bi-directional colocalization analyses between serum vitamin D and ADHD symptoms to verify their association. The results indicated that vitamin D influences ADHD symptoms on chromosome 2 (Figure 6C), whereas ADHD symptoms influence vitamin D levels on chromosome 1 (Figure 6D). In other words, serum 25-dihydroxyvitamin D3 may serve as a potential hub predictive biomarker to aid in the early diagnosis of ADHD.

Finally, we performed Gene Ontology (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway enrichment analyses to explore the potential biological mechanisms underlying the impact of vitamin D on ADHD symptoms, based on genes associated with SNPs used as instrumental variables (IVs) for serum vitamin D (Figure 6E and F). GO analysis highlighted the uronic acid metabolic process, glucuronate metabolic process, and neutral lipid metabolic process, whereas KEGG analysis highlighted the ascorbate and aldarate metabolism, pentose and glucuronate interconversions, and glycerolipid metabolism. In summary, our results suggest that serum 25-dihydroxyvitamin D3 may serve as a potential hub predictive biomarker to aid in the early diagnosis of ADHD, and that vitamin D may potentially affect ADHD symptoms via carbohydrate and lipid metabolism.

Discussion

Attention-Deficit/Hyperactivity Disorder (ADHD) is a neurodevelopmental disorder characterized by persistent patterns of inattention, hyperactivity, and impulsivity that interfere with daily functioning and development. People with ADHD may have difficulty staying on task, sustaining focus, and managing impulsive behaviors.¹² Although the number of ADHD diagnoses has sharply increased recently, underdiagnosis of ADHD is still common, which can have significant consequences, including chronic stress, low self-esteem, and difficulties in personal and professional lives.³³ Given the challenges in diagnosing ADHD accurately, identifying reliable biomarkers for ADHD is crucial. Therefore, blood-based biomarkers are being explored for their potential to aid in ADHD diagnosis, as they can help in the early detection of ADHD and lead to a better understanding of the disorder and the development of new therapeutic strategies.² Additionally, non-invasive testing simplifies the process and makes it more comfortable, especially for children with ADHD.³⁴ Recently, machine learning has been rapidly transforming the field of neurology, offering powerful tools for diagnosis, prognosis, and treatment. A previous study demonstrated that machine learning can identify individuals at high risk of developing certain neurological conditions, such as Alzheimer's disease, multiple sclerosis, and brain tumors, enabling earlier interventions and preventive strategies.³⁵ Therefore, our study aims to develop a machine learning-based risk probability predictive model for ADHD, integrating observational cohort studies with advanced machine learning algorithms to aid in early diagnosis.

In this study, we initially selected features via univariable and multivariable logistic regression analyses, which revealed that 15 features—including serum 25-dihydroxyvitamin D—were selected as predictors of ADHD risk. Then, we established a traditional predictive model using a nomogram, and the ROC, calibration, and DCA curves demonstrated good discrimination and efficacy for our feature selection. Next, we employed various machine learning algorithms to further refine the number of features and develop prediction models based on the selected clinical features from univariable and multivariable logistic regression analyses. The Gradient Boosting Machine (GBM) model demonstrated the highest performance, with an AUC of 0.91 and 0.95 in the ROC and PRC curves respectively. Eight features, including serum 25-dihydroxyvitamin D, were identified as the optimal set for the GBM model. The Shapley Additive eXplanation (SHAP) method was employed to interpret this model, revealing

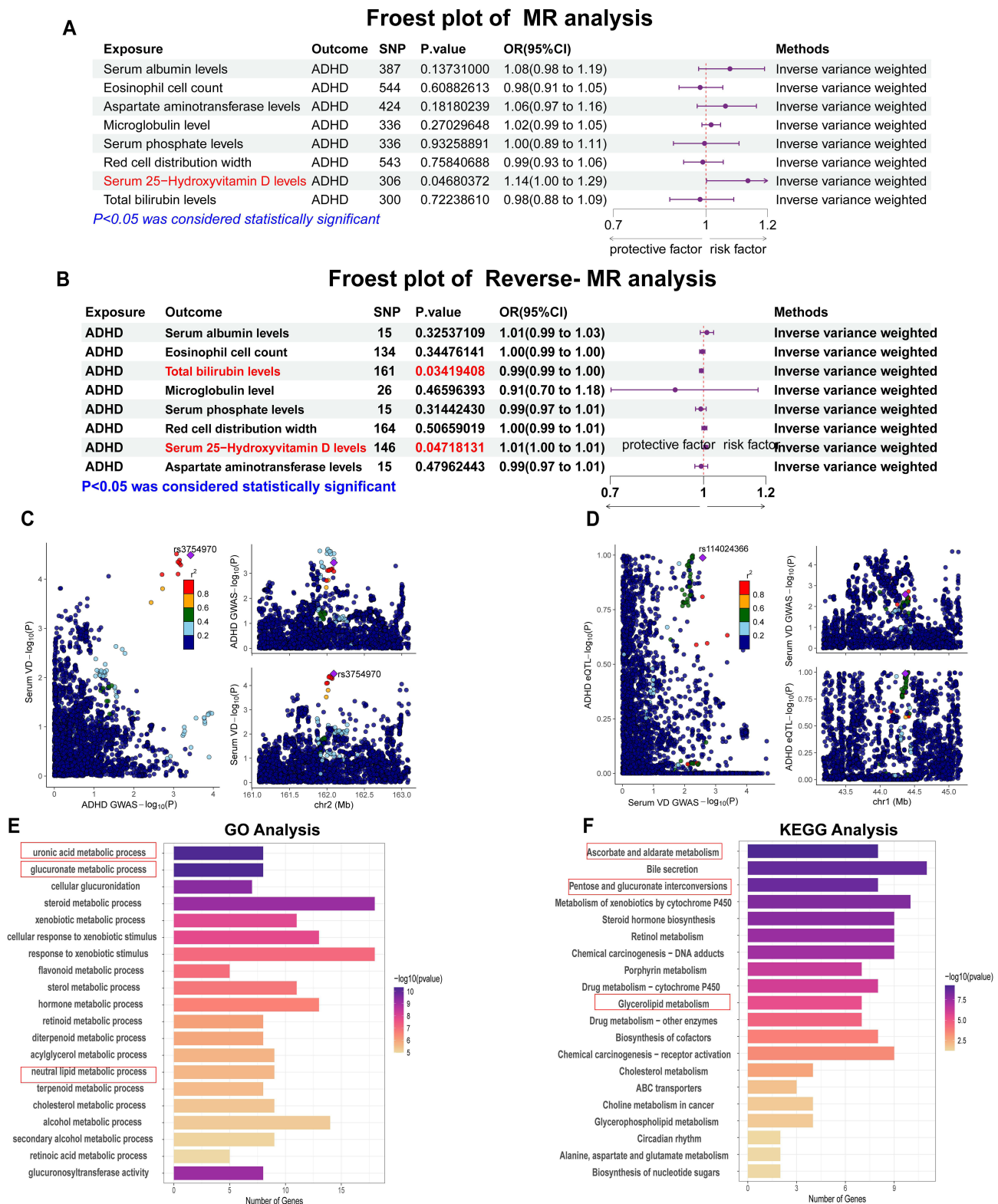


Figure 6 Bi-directional Mendelian Randomization analysis. **(A and B)** Forest plot of Mendelian randomization and reverse Mendelian randomization analysis. **(C and D)** Bi-directional colocalization analysis between serum vitamin D and ADHD symptoms, suggesting that vitamin D can serve as a biomarker of ADHD. **(E and F)** Gene Ontology (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway enrichment analysis revealed significant pathways involving genes associated with SNPs using as the instrumental variables (IVs) of serum vitamin D.

the pivotal role of the top four features—blood beta-2 microglobulin levels, red blood cell distribution width, 25-dihydroxyvitamin D3, and the percentage of eosinophils—in predicting ADHD risk probability. The final prediction model was deployed as a web application to facilitate early diagnosis and intervention in clinical scenarios. However, it should be used as a supportive tool alongside comprehensive clinical assessments, and additional validation remains necessary.

Bi-directional Mendelian randomization analyses were conducted to explore and verify the association between the selected eight features from the GBM prediction model and ADHD symptoms, based on our previous study. The results indicated that serum 25-dihydroxyvitamin D3 was associated with ADHD symptoms in both directions, and that ADHD symptoms may affect total bilirubin levels. Recently, increasing evidence has supported the crucial role of vitamin D in various neurological disorders,³⁶ and research indicates that vitamin D may play a role in protecting against Parkinson's disease by supporting neuronal health and reducing inflammation.³⁶ Moreover, there is growing evidence suggesting a link between serum vitamin D levels and ADHD. Studies have consistently found that children and adolescents with ADHD tend to have significantly lower serum concentrations of 25-dihydroxyvitamin D3 compared to healthy controls.³⁷ Although there is still limited research directly linking total bilirubin levels to ADHD, elevated bilirubin levels can influence neuroinflammation and trigger neuroinflammatory responses.³⁸ This neuroinflammation might influence neurological conditions, including ADHD.³⁹

The bi-directional colocalization analyses also supported the association of serum 25-dihydroxyvitamin D3 with ADHD symptoms in both directions. These results highlight the potential value of serum 25-dihydroxyvitamin D3 as a biomarker for predicting the risk probability of ADHD. Gene Ontology (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway enrichment analyses were conducted to explore the potential biological mechanisms underlying the impact of vitamin D on ADHD symptoms based on genes associated with SNPs used as instrumental variables (IVs) for serum vitamin D. GO analysis highlighted the uronic acid metabolic process, glucuronate metabolic process, and neutral lipid metabolic process, and KEGG analysis highlighted the ascorbate and aldarate metabolism, pentose and glucuronate interconversions, and glycerolipid metabolism. Research has indicated that children with ADHD may have altered carbohydrate metabolism, including changes in the levels of certain uronic acids.⁴⁰ Uronic acids, such as glucuronic acid, are involved in detoxifying various substances in the body, particularly certain environmental pollutants like Bisphenol-A (BPA) and Diethylhexyl Phthalate (DEHP). A recent study involving children with ADHD, autism spectrum disorder (ASD), and neurotypical controls indicated that the efficiency of glucuronidation for BPA was reduced by about 17% in children with ADHD compared to the control group.⁴¹ Glycerolipid metabolism involves the synthesis and breakdown of glycerolipids, which include triglycerides and phospholipids.⁴² This process is crucial for maintaining cellular energy balance and producing signaling molecules.⁴² Neutral lipids, such as triglycerides, play a significant role in this metabolic pathway.⁴³ In the context of ADHD, research has shown that individuals with ADHD may have imbalances in their lipid profiles, such as higher levels of triglycerides and lower levels of HDL (good cholesterol).^{44,45} These imbalances can influence brain function and behavior, highlighting the importance of lipid metabolism in the context of ADHD. Ascorbate and aldarate metabolism are crucial pathways for maintaining cellular health and protecting against oxidative damage.⁴⁶ In the context of ADHD, antioxidants play a significant role in mitigating oxidative stress, which is often elevated in individuals with ADHD.⁴⁷ Oxidative stress can damage cells and tissues, including those in the brain, potentially exacerbating ADHD symptoms. Additionally, recent studies indicate that antioxidant therapy might improve symptoms in both ADHD and epilepsy (potentially co-occurring disorders) by reducing oxidative damage and inflammation.^{48,49} Furthermore, advanced studies have indicated that individuals with ADHD may experience difficulties with theory of mind skills—the cognitive ability to understand that others have their own mental states. Children with ADHD also exhibit significantly altered brain activity compared to typically developing controls, including increased amplitude of low-frequency fluctuation and decreased functional connectivity in various brain regions.^{50,51} Growing evidence suggests that oxidative stress may contribute to the pathophysiology of these symptoms, and vitamin D may alleviate them due to its antioxidant and anti-inflammatory properties within the brain.⁵² Taken together, these findings suggest that serum 25-dihydroxyvitamin D3 may serve as a central biomarker for ADHD risk, with vitamin D potentially influencing ADHD symptoms via its roles in carbohydrate metabolism, lipid metabolism, and regulation of the antioxidant system.

Limitation

There are several limitations to this study. First, this study was conducted at a single center in China and involved Chinese populations, so the findings may not be directly applicable to other populations or ethnicities. Second, this study developed a machine learning-based prediction model with 740 participants, which is a decent sample size. Further external validation is still required to assess the generalizability of this model. Third, the pilot prospective observational cohort study covered a relatively short period. Further longitudinal data and follow-up assessments are required to provide more robust evidence for the utility of serum 25-dihydroxyvitamin D3 as a biomarker for ADHD. Fourth, ADHD is a complex neurodevelopmental disorder, and a single biomarker may not be sufficient to capture its full complexity. Further research investigating other potential biomarkers could provide a more comprehensive understanding of ADHD pathogenesis. Fifth, this study does not provide direct mechanistic insights into how vitamin D deficiency contributes to ADHD pathogenesis. Further experimental studies are needed to elucidate the underlying biological mechanisms. Finally, although this study suggests that serum 25-dihydroxyvitamin D3 could be a promising biomarker for ADHD, its clinical utility and cost-effectiveness need to be further evaluated in real-world settings before it can be widely implemented in clinical practice.

Conclusion

In conclusion, we successfully developed a machine learning (ML) model to predict the risk probability of ADHD using clinical data obtained from real-world clinical practice. The GBM model demonstrated superior performance compared with five other machine learning algorithms in this study. Additionally, the Shapley Additive eXplanation (SHAP) method was applied to elucidate the ML model. This approach not only helped determine the importance of each feature in the model but also demonstrated how each feature influenced the model. Finally, we explored and verified the association between the selected eight features from the ML model and ADHD symptoms using bi-directional Mendelian randomization analyses. Taken together, our study successfully establishes a machine learning-based risk probability predictive model for ADHD, indicating that serum 25-dihydroxyvitamin D3 may serve as a potential hub predictive biomarker to aid in the disorder's early diagnosis and that vitamin D may potentially affect ADHD symptoms through carbohydrate and lipid metabolism pathways.

Data Sharing Statement

The raw code used in this study is openly available in our GitHub repository: <https://github.com/PediatricLab-YueyingLiu/ADHDML>. Additional data will be made available upon reasonable request to the corresponding author.

Acknowledgments

We would like to express our sincere gratitude to all the children with ADHD who generously contributed their time and effort to our pilot cohort. Their participation and dedication have been instrumental to the successful completion of this research. We also extend our heartfelt thanks to their parents and guardians for their trust and support throughout the study. We deeply appreciate the medical staff and research assistants who played vital roles in the execution of this project. Their expertise and commitment were invaluable to our research process. This study was supported by a grant from the National Natural Science Foundation of China awarded to Yueying Liu (No. 82371462), the Qing Lan Project of Jiangsu to Ruijin Xie (JS2023-27), and the Traditional Chinese Medicine Science and Technology Development Plan Project of Zhejiang Province to Yuanying Jiang (Grant No. 2024ZL140). We gratefully acknowledge this financial support, which made our research possible. We would like to thank the School of Medicine at Jiangnan University for providing the facilities and resources necessary to conduct this study.

Disclosure

All authors confirm that there are no conflicts of interest in this work.

References

- Salari N, Ghasemi H, Abdoli N, et al. The global prevalence of ADHD in children and adolescents: a systematic review and meta-analysis. *Ital J Pediatr.* 2023;49(1):48. doi:10.1186/s13052-023-01456-1
- Michellini G, Norman LJ, Shaw P, Loo SK. Treatment biomarkers for ADHD: taking stock and moving forward. *Transl Psychiatry.* 2022;12(1):444. doi:10.1038/s41398-022-02207-2
- Cynthia Reuben MA, Nazik Elgaddal MS. Attention-deficit/hyperactivity disorder in children ages 5–17 Years: United States, 2020–2022. Available from: <https://www.cdc.gov/nchs/products/databriefs/db499.htm>. Accessed May 09, 2025.
- Gomez R, Chen W, Houghton S. Differences between DSM-5-TR and ICD-11 revisions of attention deficit/hyperactivity disorder: a commentary on implications and opportunities. *World J Psychiatry.* 2023;13(5):138–143. doi:10.5498/wjp.v13.i5.138
- de Lima TA, Zuanetti PA, Nunes MEN, Hamad APA. Differential diagnosis between autism spectrum disorder and other developmental disorders with emphasis on the preschool period. *World J Pediatr.* 2023;19(8):715–726. doi:10.1007/s12519-022-00629-y
- Al-Yateem N, Slewa-Younan S, Halimi A, et al. Prevalence of undiagnosed attention deficit hyperactivity disorder (ADHD) symptoms in the young adult population of the United Arab Emirates: a national cross-sectional study. *J Epidemiol Glob Health.* 2024;14(1):45–53. doi:10.1007/s44197-023-00167-4
- Martínez-Jaime MM, Reyes-Morales H, Peyrot-Negrete I, Barrientos-álvarez MS. Access to early diagnosis for attention-deficit/hyperactivity disorder among children and adolescents in Mexico City at specialized mental health services. *BMC Health Serv Res.* 2024;24(1):599. doi:10.1186/s12913-024-11022-y
- Jepsen IB, Brynskov C, Thomsen PH, Rask CU, Jensen de López K, Lambek R. The role of language in the social and academic functioning of children with ADHD. *J Atten Disord.* 2024;28(12):1542–1554. doi:10.1177/10870547241266419
- Wei Q, Xu X, Xu X, Cheng Q. Early identification of autism spectrum disorder by multi-instrument fusion: a clinically applicable machine learning approach. *Psychiatry Res.* 2023;320:115050. doi:10.1016/j.psychres.2023.115050
- Mirhosseini H, Maayeshi N, Hooshmandi H, Moradkhani S, Hosseinzadeh M. The effect of vitamin D supplementation on the brain mapping and behavioral performance of children with ADHD: a double-blinded randomized controlled trials. *Nutr Neurosci.* 2024;27(6):566–576. doi:10.1080/1028415x.2023.2233752
- Mei H, Xie R, Li T, Chen Z, Liu Y, Sun C. Effect of atomoxetine on behavioral difficulties and growth development of primary school children with attention-deficit/hyperactivity disorder: a prospective study. *Children.* 2022;9(2). doi:10.3390/children9020212
- Koutsoklenis A, Honkasilta J. ADHD in the DSM-5-TR: what has changed and what has not. *Front Psychiatry.* 2022;13:1064141. doi:10.3389/fpsy.2022.1064141
- Faraone SV, Bellgrove MA, Brikell I, et al. Attention-deficit/hyperactivity disorder. *Nat Rev Dis Primers.* 2024;10(1):11. doi:10.1038/s41572-024-00495-0
- Hu J, Xu J, Li M, et al. Identification and validation of an explainable prediction model of acute kidney injury with prognostic implications in critically ill children: a prospective multicenter cohort study. *EClinicalMedicine.* 2024;68:102409. doi:10.1016/j.eclinm.2023.102409
- Guo X, Ke Y, Wu B, et al. Exploratory analysis of the association between organophosphate ester mixtures with high blood pressure of children and adolescents aged 8–17 years: cross-sectional findings from the national health and nutrition examination survey. *Environ Sci Pollut Res Int.* 2023;30(9):22900–22912. doi:10.1007/s11356-022-23740-z
- Võsa U, Claringbould A, Westra HJ, et al. Large-scale cis- and trans-eQTL analyses identify thousands of genetic loci and polygenic scores that regulate blood gene expression. *Nat Genet.* 2021;53(9):1300–1310. doi:10.1038/s41588-021-00913-z
- Teng Z, Li L, Xin Z, et al. A literature review of artificial intelligence (AI) for medical image segmentation: from AI and explainable AI to trustworthy AI. *Quant Imaging Med Surg.* 2024;14(12):9620–9652. doi:10.21037/qims-24-723
- Ning Y, Ong MEH, Chakraborty B, et al. Shapley variable importance cloud for interpretable machine learning. *Patterns.* 2022;3(4):100452. doi:10.1016/j.patter.2022.100452
- Nápoles-Duarte JM, Biswas A, Parker MI, Palomares-Baez JP, Chávez-Rojó MA, Rodríguez-Valdez LM. Stmol: a component for building interactive molecular visualizations within streamlit web-applications. *Front Mol Biosci.* 2022;9:990846. doi:10.3389/fmolb.2022.990846
- Sanderson E, Glymour MM, Holmes MV, et al. Mendelian randomization. *Nat Rev Meth Primers.* 2022;2(1). doi:10.1038/s43586-021-00092-5
- Lin L, Zhang R, Huang H, et al. Mendelian randomization with refined instrumental variables from genetic score improves accuracy and reduces Bias. *Front Genet.* 2021;12:618829. doi:10.3389/fgene.2021.618829
- Liu Y, Chang Y, Jiang X, et al. Analysis of the role of PANoptosis in seizures via integrated bioinformatics analysis and experimental validation. *Heliyon.* 2024;10(4):e26219. doi:10.1016/j.heliyon.2024.e26219
- Xie Q, Hu B. Effects of gut microbiota on prostatic cancer: a two-sample Mendelian randomization study. *Front Microbiol.* 2023;14:1250369. doi:10.3389/fmicb.2023.1250369
- Cao Y, Zhao W, Zhong Y, et al. Effects of chronic low-level lead (Pb) exposure on cognitive function and hippocampal neuronal ferroptosis: an integrative approach using bioinformatics analysis, machine learning, and experimental validation. *Sci Total Environ.* 2024;917:170317. doi:10.1016/j.scitotenv.2024.170317
- Lopera-Maya EA, Kurilshikov A, van der Graaf A, et al. Effect of host genetics on the gut microbiome in 7,738 participants of the Dutch microbiome project. *Nat Genet.* 2022;54(2):143–151. doi:10.1038/s41588-021-00992-y
- Oscanio J, Sivapalan L, Gadaleta E, Dayem Ullah AZ, Lemoine NR, Chelala C. SNPnexus: a web server for functional annotation of human genome sequence variation (2020 update). *Nucleic Acids Res.* 2020;48(W1):W185–w192. doi:10.1093/nar/gkaa420
- Xu S, Hu E, Cai Y, et al. Using clusterProfiler to characterize multiomics data. *Nat Protoc.* 2024;19(11):3292–3320. doi:10.1038/s41596-024-01020-z
- Yu G. Thirteen years of clusterProfiler. *Innovation.* 2024;5(6):100722. doi:10.1016/j.xinn.2024.100722
- Mei H, Wu D, Yong Z, et al. PM(2.5) exposure exacerbates seizure symptoms and cognitive dysfunction by disrupting iron metabolism and the Nrf2-mediated ferroptosis pathway. *Sci Total Environ.* 2024;910:168578. doi:10.1016/j.scitotenv.2023.168578
- Michaëlsson M, Yuan S, Melhus H, et al. The impact and causal directions for the associations between diagnosis of ADHD, socioeconomic status, and intelligence by use of a bi-directional two-sample Mendelian randomization design. *BMC Med.* 2022;20(1):106. doi:10.1186/s12916-022-02314-3

31. Keilow M, Holm A, Fallesen P. Medical treatment of attention deficit/hyperactivity disorder (ADHD) and children's academic performance. *PLoS One*. 2018;13(11):e0207905. doi:10.1371/journal.pone.0207905
32. Arnold IC, Munitz A. Spatial adaptation of eosinophils and their emerging roles in homeostasis, infection and disease. *Nat Rev Immunol*. 2024;24(12):858–877. doi:10.1038/s41577-024-01048-y
33. French B, Daley D, Groom M, Cassidy S. Risks associated with undiagnosed ADHD and/or autism: a mixed-method systematic review. *J Atten Disord*. 2023;27(12):1393–1410. doi:10.1177/10870547231176862
34. Takahashi N, Ishizuka K, Inada T. Peripheral biomarkers of attention-deficit hyperactivity disorder: current status and future perspective. *J Psychiatr Res*. 2021;137:465–470. doi:10.1016/j.jpsychires.2021.03.012
35. Kalani M, Anjankar A. Revolutionizing neurology: the role of artificial intelligence in advancing diagnosis and treatment. *Cureus*. 2024;16(6):e61706. doi:10.7759/cureus.61706
36. Plantone D, Primiano G, Manco C, Locci S, Servidei S, De Stefano N. Vitamin D in Neurological Diseases. *Int J Mol Sci*. 2022;24(1). doi:10.3390/ijms24010087
37. Lukovac T, Hil OA, Popović M, et al. Serum biomarker analysis in pediatric ADHD: implications of homocysteine, vitamin B12, vitamin D, ferritin, and iron levels. *Children*. 2024;11(4). doi:10.3390/children11040497
38. Zhang F, Chen L, Jiang K. Neuroinflammation in Bilirubin neurotoxicity. *J Integr Neurosci*. 2023;22(1):9. doi:10.31083/j.jin2201009
39. Jayanti S, Dalla Verde C, Tiribelli C, Gazzin S. Inflammation, dopaminergic brain and Bilirubin. *Int J Mol Sci*. 2023;24(14). doi:10.3390/ijms241411478
40. Roetner J, Van Doren J, Maschke J, et al. Effects of prenatal alcohol exposition on cognitive outcomes in childhood and youth: a longitudinal analysis based on meconium ethyl glucuronide. *Eur Arch Psychiatry Clin Neurosci*. 2024;274(2):343–352. doi:10.1007/s00406-023-01657-z
41. Stein TP, Schluter MD, Steer RA, Ming X. Bisphenol-A and phthalate metabolism in children with neurodevelopmental disorders. *PLoS One*. 2023;18(9):e0289841. doi:10.1371/journal.pone.0289841
42. Watanabe Y, Kasuga K, Tokutake T, Kitamura K, Ikeuchi T, Nakamura K. Alterations in glycerolipid and fatty acid metabolic pathways in Alzheimer's disease identified by urinary metabolic profiling: a pilot study. *Front Neurol*. 2021;12:719159. doi:10.3389/fneur.2021.719159
43. Farese RV Jr, Walther TC. Glycerolipid synthesis and lipid droplet formation in the endoplasmic reticulum. *Cold Spring Harb Perspect Biol*. 2023;15(5):a041246. doi:10.1101/cshperspect.a041246
44. Chen X, Zheng Z, Xie D, et al. Serum lipid metabolism characteristics and potential biomarkers in patients with unilateral sudden sensorineural hearing loss. *Lipids Health Dis*. 2024;23(1):205. doi:10.1186/s12944-024-02189-8
45. Yang D, Wang X, Zhang L, et al. Lipid metabolism and storage in neuroglia: role in brain development and neurodegenerative diseases. *Cell Biosci*. 2022;12(1):106. doi:10.1186/s13578-022-00828-0
46. Liang H, Song K. Elucidating ascorbate and aldarate metabolism pathway characteristics via integration of untargeted metabolomics and transcriptomics of the kidney of high-fat diet-fed obese mice. *PLoS One*. 2024;19(4):e0300705. doi:10.1371/journal.pone.0300705
47. Corona JC. Role of oxidative stress and neuroinflammation in attention-deficit/hyperactivity disorder. *Antioxidants*. 2020;9(11). doi:10.3390/antiox9111039
48. Zhou P, Yu X, Song T, Hou X. Safety and efficacy of antioxidant therapy in children and adolescents with attention deficit hyperactivity disorder: a systematic review and network meta-analysis. *PLoS One*. 2024;19(3):e0296926. doi:10.1371/journal.pone.0296926
49. de Melo AD, Freire VAF, Diogo ÍL, Santos HL, Barbosa LA, de Carvalho LED. Antioxidant therapy reduces oxidative stress, restores Na, K-ATPase function and induces neuroprotection in rodent models of seizure and epilepsy: a systematic review and meta-analysis. *Antioxidants*. 2023;12(7). doi:10.3390/antiox12071397
50. Kılınçel Ş. The relationship between the theory of mind skills and disorder severity among adolescents with ADHD. *Alpha Psychiatry*. 2021;22(1):7–11. doi:10.5455/apd.126537
51. Liao W, Li H, Liu Q, et al. Comparison of brain function between medication-naïve adhd with and without comorbidity in Chinese children using resting-state fNIRS. *Alpha Psychiatry*. 2024;25(4):485–492. doi:10.5152/alphapsychiatry.2024.241674
52. Renke G, Starling-Soares B, Baesso T, Petronio R, Aguiar D, Paes R. Effects of vitamin D on cardiovascular risk and oxidative stress. *Nutrients*. 2023;15(3):769. doi:10.3390/nu15030769

Neuropsychiatric Disease and Treatment

Publish your work in this journal

Neuropsychiatric Disease and Treatment is an international, peer-reviewed journal of clinical therapeutics and pharmacology focusing on concise rapid reporting of clinical or pre-clinical studies on a range of neuropsychiatric and neurological disorders. This journal is indexed on PubMed Central, the 'PsycINFO' database and CAS, and is the official journal of The International Neuropsychiatric Association (INA). The manuscript management system is completely online and includes a very quick and fair peer-review system, which is all easy to use. Visit <http://www.dovepress.com/testimonials.php> to read real quotes from published authors.

Submit your manuscript here: <https://www.dovepress.com/neuropsychiatric-disease-and-treatment-journal>

Dovepress
Taylor & Francis Group