

# Project Overview: Mental Health Treatment Prediction Chatbot

## Project Goal

The primary objective of this project was to develop a highly accurate and reliable machine learning model to predict the likelihood of an individual seeking professional mental health treatment. This model serves as the core intelligence for a conversational chatbot designed to provide personalized recommendations and guidance based on a user's self-reported profile and concerns.

The core success criteria were twofold:

1. **High Predictive Performance:** Achieve a classification model with strong discriminatory power (high ROC-AUC and F1-Score).
2. **Model Interpretability:** Ensure the final model allows for clear insight into the most influential factors, enabling the chatbot to offer explainable and trustworthy recommendations.

## Final Model and Performance

The final production model selected is the **XGBoost Classifier**, thoroughly optimized using the **Optuna** hyperparameter tuning framework on a dataset enhanced with **Target Encoding**.

The model demonstrated exceptional performance on the validation set:

ROC Curve: 0.9858    Outstanding ability to distinguish between the two classes (treatment vs. no treatment).

Accuracy: 0.9200    Correctly classifies individuals 92% of the time.

F1 Score: 0.9202    Excellent balance between false positives and false negatives, crucial for sensitive medical recommendations.

This high performance confirms the model's robustness and readiness for deployment within the chatbot application.

## Challenges Faced and Data Sourcing Process

The project's development phase was marked by significant challenges, particularly in securing a suitable, high-quality dataset, a common hurdle in real-world ML projects.

### *1. Data Selection Hurdles (The Four Failed Attempts)*

Finding a dataset that was both relevant and analytically viable required extensive searching and validation:

The rejected data sets and the cause of rejection:

**Brain Stroke: Extreme Imbalance (5% Class 1):** The overwhelming class imbalance made it difficult to train a model that could generalize well to the minority class without aggressive, and often unreliable, sampling techniques.

**Germany Heart Disease: Synthetic Data:** The dataset consisted of simulated, non-real data, leading to concerns about the authenticity and generalizability of any insights derived, making it unsuitable for a serious application.

**Rejected Mental Illness: Missing Target Variable:** The dataset lacked a crucial column defining the outcome or target (e.g., whether treatment was sought), making supervised learning (model training) impossible.

**Appointment No Show: Significant Imbalance (29% Class 1):** While less severe than the first attempt, the 29% imbalance required substantial and ultimately unsuccessful efforts with multiple balancing techniques (e.g., SMOTE, undersampling) that failed to produce stable, high-performing metrics.

## *2. The Solution: The Mental Health Dataset*

The project moved forward successfully with the **Mental Health Data set**, which provided a suitable balance of features and a tractable class distribution. This dataset served as the foundation for the final model.

## *3. Technical Challenges*

- **Handling High-Cardinality Features (e.g., Country):** Features with a large number of unique values, such as the user's `Country`, posed a challenge for standard encoding techniques (like One-Hot Encoding), which would have created hundreds of sparse, low-importance columns. This was overcome by employing **Target Encoding**, which summarized the target likelihood (treatment rate) for each country, creating one highly informative numeric feature.
- **Model Optimization:** Given the complexity of the data, achieving the high benchmark metrics required rigorous optimization. This was solved by integrating the **Optuna** framework to systematically search and discover the globally best hyperparameters for the XGBoost model, moving beyond standard grid or random search.

# Data Acquisition, Preprocessing, and Initial Analysis

## 1. Introduction

This section details the initial data handling, cleaning, and exploratory analysis performed on the **Mental Health Dataset**. The objective was to prepare the raw data for modeling by identifying inconsistencies, removing redundancies, and understanding the underlying distributions of key variables.

## 2. Data Loading and Structural Inspection

The dataset was loaded using the Pandas library. An initial inspection of the dataframe structure revealed the following:

- **Initial Dimensions:** The raw dataset contained **292,364** entries and **17** columns.
- **Data Types:** The dataset consists primarily of categorical data (`object` type), including demographic details (Gender, Country) and survey responses (e.g., `Growing_Stress`, `Work_Interest`).

## 3. Data Cleaning and Integrity Checks

### 3.1 Duplicate Removal

To ensure the statistical validity of the analysis, the dataset was screened for duplicate entries.

- **Duplicates Identified:** A total of **2,313** duplicate rows were detected.
- **Action Taken:** These duplicates were removed to prevent data leakage and bias.
- **Final Dimensions:** The cleaned dataset contains **290,051** unique records.

### 3.2 Missing Value Analysis

A null value assessment was conducted across all columns.

- **Findings:** The `self_employed` column was identified as having **5,193** missing values.
- **Integrity:** All other columns (`Timestamp`, `Gender`, `Country`, etc.) were found to be complete with zero null values.

## 4. Feature Analysis

To facilitate analysis, features were categorized based on their data types and cardinality.

- **Categorical Identification:** Columns were automatically flagged as categorical if they were of the `object` data type or had fewer than 10 unique values.
- **Value Counts:** Frequency tables were generated for each categorical column to check for cardinality and potential data entry errors (e.g., checking for inconsistent naming conventions in the `Gender` column).

## 5. Univariate Analysis (Visualizations)

A comprehensive univariate analysis was performed to visualize the distribution of each categorical variable.

- **Methodology:** Iterative plotting was used to generate histograms for all key demographic and mental health indicators, including `Occupation`, `Days_Indoors`, and `Growing_Stress`.
- **Visualization Tool:** Plotly Express was utilized to create interactive charts, allowing for detailed inspection of specific categories.
- **Goal:** These visualizations serve to identify class imbalances (e.g., if one country is overrepresented) and general trends in the survey responses before moving to multivariate analysis.

## 6. Feature Engineering and Transformation

Following the initial cleaning, the dataset underwent extensive feature engineering to transform raw variables into a machine-learning-ready format.

### 6.1 Temporal Feature Extraction

The `Timestamp` column was decomposed into constituent numerical components: `Year`, `Month`, `Day`, and `Hour`. To capture potential seasonal or behavioral patterns, three boolean features were derived:

- **Is\_Winter:** Flagging months 12, 1, and 2 (assuming Northern Hemisphere winters).
- **Is\_MidYear:** Flagging the period from May to August.
- **Is\_Night:** Flagging hours between 20:00 and 06:00.

### 6.2 Categorical Variable Encoding

Categorical variables were mapped to numerical representations based on their ordinality and cardinality:

- **Binary Mapping:** Variables such as `Gender`, `family_history`, `treatment`, and `Coping_Struggles` were mapped to binary values (0/1).
  - *Note on Missing Values:* Missing values in `self_employed` were imputed as 0 (No) before mapping, under the assumption that non-responses likely indicate traditional employment.
- **Ordinal Mapping:** `Days_Indoors` and `Mood_Swings` were mapped to integer scales to preserve the inherent order of their categories (e.g., '1-14 days' → 1, 'More than 2 months' → 4).
- **Label Encoding:** `Growing_Stress`, `Changes_Habits`, `Mental_Health_History`, `Work_Interest`, `Social_Weakness`, `mental_health_interview`, and `care_options` were transformed using Label Encoding.
- **One-Hot Encoding:** The `Occupation` variable was One-Hot Encoded to handle its nominal nature without imposing an artificial order.

### 6.3 Target Encoding with Cross-Validation

The `Country` variable, which possesses high cardinality, was encoded using the mean target value (`treatment rate`).

- **Leakage Prevention:** To prevent data leakage, a 5-fold Stratified Cross-Validation strategy was employed on the training set. The mean treatment rate was calculated for each fold and mapped to the corresponding validation fold.
- **Test Set Application:** The global mean treatment rates from the training set were then mapped to the test set, with unseen countries filled using the global mean.

## 6.4 Composite Feature Creation

Several interaction features were engineered to capture complex relationships:

- **Stress\_Score:** An aggregate mean of five stress-related indicators (including `Days_Indoors` and `Mood_Swings`).
- **Social\_Function\_Score:** Calculated as `Work_Interest` minus `Social_Weakness`.
- **SelfEmployment\_Risk:** An interaction term combining `self_employed` status with the lack of `care_options`.
- **Family\_Support\_Impact:** An interaction term combining `family_history` with `Coping_Struggles`.

## 7. Data Splitting and Final Output

To prepare for model training, the processed dataframe was split into training and testing sets.

- **Split Ratio:** 80% Training, 20% Testing.
- **Stratification:** The split was stratified based on the `treatment` target variable to ensure class distributions remained consistent across both sets.
- **Final Artifacts:** The resulting datasets were exported as `train_encoded.csv` and `test_encoded.csv`.

## Baseline Model Selection and Evaluation

### Model 1: CatBoost Classifier (Baseline)

The CatBoost Classifier, a gradient boosting decision tree algorithm, was selected as the primary strong baseline model due to its robustness against overfitting and its handling of various feature types. The model was trained on the target-encoded feature set with optimized hyperparameters (300 iterations, depth of 8, learning rate of 0.05).

*Performance Metrics on Test Set:*

ROC Curve: 0.9844

Accuracy: 0.9191

Precision: 0.92508

Recall: 0.9135

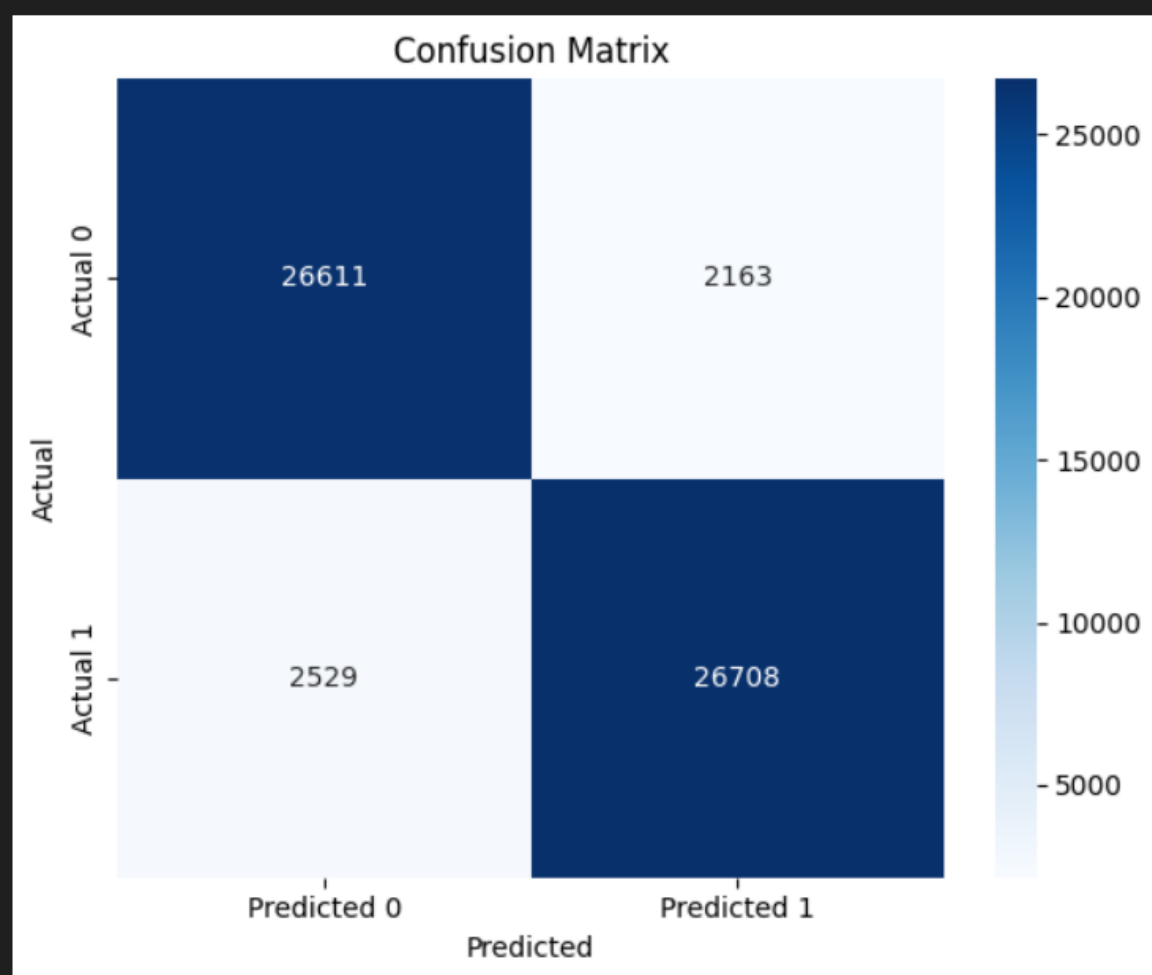
F1 Score: 0.9193

The classification report provides a detailed breakdown:

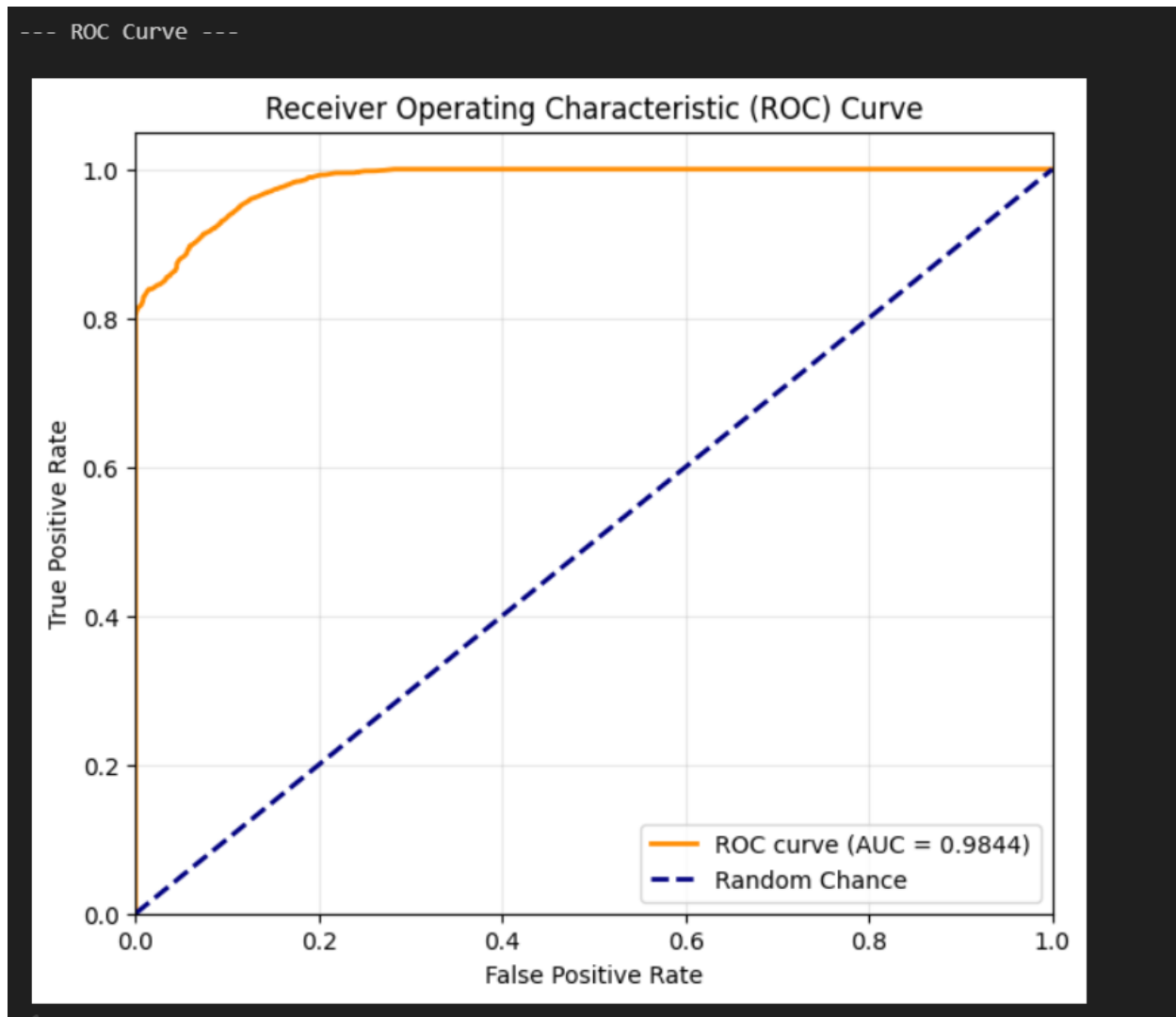
```
--- Classification Report ---
              precision    recall  f1-score   support

     0       0.9132       0.9248       0.9190       28774
     1       0.9251       0.9135       0.9193       29237

 accuracy          0.9191          0.9191       58011
 macro avg       0.9191       0.9192       0.9191       58011
weighted avg       0.9192       0.9191       0.9191       58011
```



### *ROC Curve*



### *Feature Importance Analysis*

The model provided strong insights into the most influential features driving the prediction of treatment need.



```
--- Feature Importance ---
Top 20 most important features:
      Feature Id    Import
      Hour  24.121832
      care_options 16.464995
Country_TreatmentRate 14.623283
      Day 13.973181
      family_history 12.153335
mental_health_interview 6.993640
      Gender 5.281923
      self_employed 2.262002
      Is_Night 1.997255
SelfEmployment_Risk 1.039574
      Is_MidYear 0.498296
      Year 0.427103
      Month 0.133551
      Is_Winter 0.024737
Family_Support_Impact 0.001472
      Stress_Score 0.000963
      Work_Interest 0.000915
      Changes_Habits 0.000511
Mental_Health_History 0.000391
      Days_Indoors 0.000363

--- Evaluation Complete ---
```

## Model 2: CatBoost Classifier (Raw Categorical Features)

To establish a benchmark for the effectiveness of manual feature engineering and to mitigate potential data leakage from target encoding, a second CatBoost model was trained directly on the raw, unencoded categorical features. CatBoost's native feature handling was leveraged, requiring only time-series decomposition.

### *Performance Metrics on Test Set*

ROC Curve: 0.9868

Accuracy: 0.9223

Precision: 0.9174

Recall: 0.9284

F1 Score: 0.9229

The classification report provides a detailed breakdown:

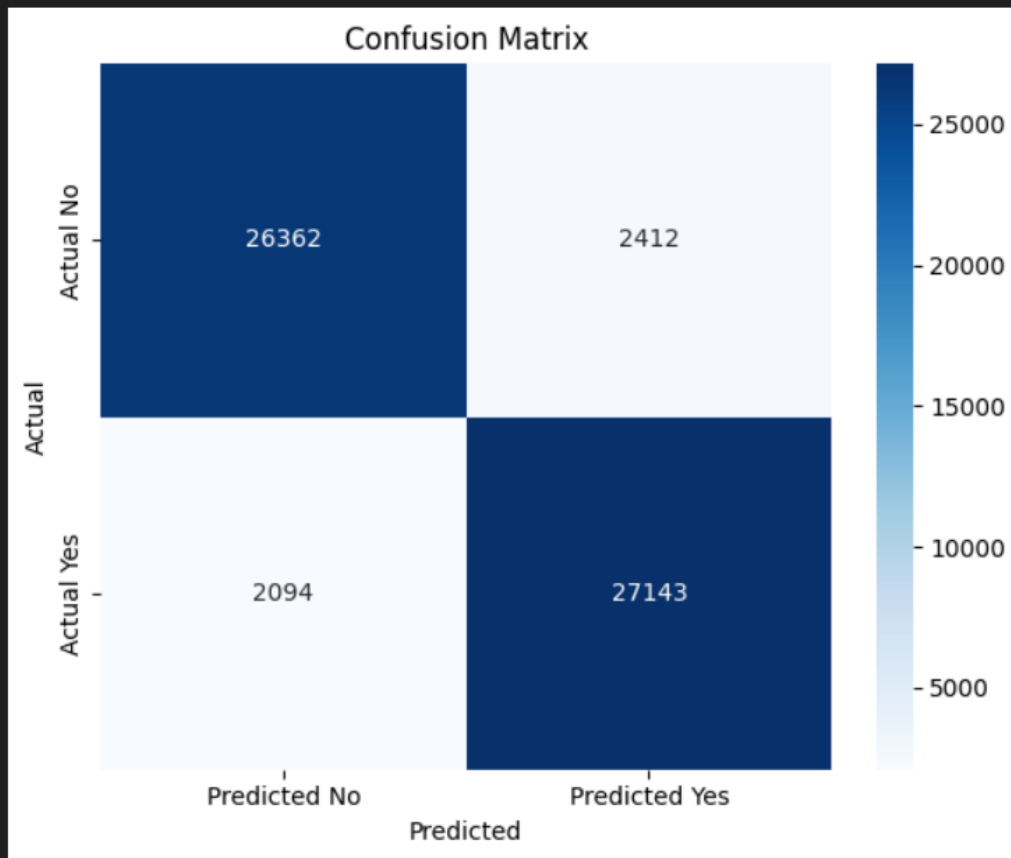
```
--- Classification Report ---
              precision    recall  f1-score   support

     0       0.9264       0.9162       0.9213       28774
     1       0.9184       0.9284       0.9234       29237

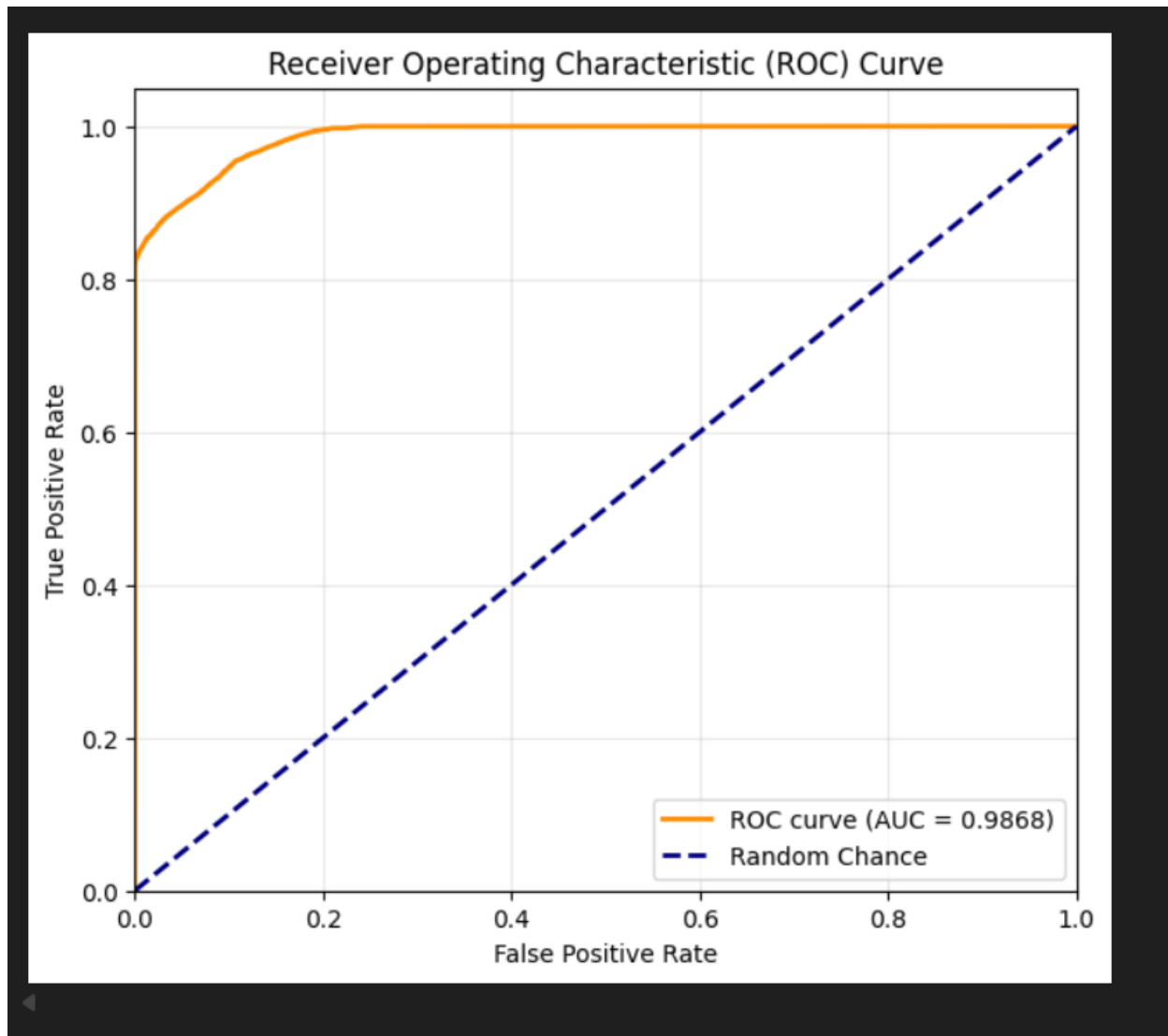
 accuracy          0.9223       58011
 macro avg       0.9224       0.9223       0.9223       58011
weighted avg       0.9224       0.9223       0.9223       58011
```

--- Confusion Matrix ---

Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings...](#)



## ROC Curve



### *Feature Importance Analysis*

```
--- Feature Importance ---
Top 20 most important features:
      Feature Id  Importances
      ts_hour    22.870466
      care_options 21.018034
      Country     12.429106
      family_history 11.629695
      ts_dayofweek 10.002504
mental_health_interview 7.240805
      Gender      7.039018
      self_employed 6.002084
      ts_month     1.542882
      Growing_Stress 0.069542
      Changes_Habits 0.054335
      Mental_Health_History 0.039505
      Occupation   0.027329
      Social_Weakness 0.016987
      Days_Indoors  0.009511
      Mood_Swings   0.005578
      Work_Interest 0.002360
      Coping_Struggles 0.000260

--- Evaluation Complete ---
```

### **Model 3: Stacking Classifier (Modified SKLearn Ensemble)**

Due to missing dependency issues with CatBoost, a modified Stacking Classifier was trained using only standard scikit-learn models that can utilize the pre-scaled, target-encoded dataset from Model 1's preparation pipeline. This ensemble combines the non-linear decision boundary capabilities of a **Random Forest** and a **Multi-Layer Perceptron (MLP)**, which are fed into a **Logistic Regression** meta-model.

### *Performance Metrics on Test Set*

ROC Curve: 0.9823

Accuracy: 0.9146

Precision: 0.9023

Recall: 0.9313

F1 Score: 0.9166

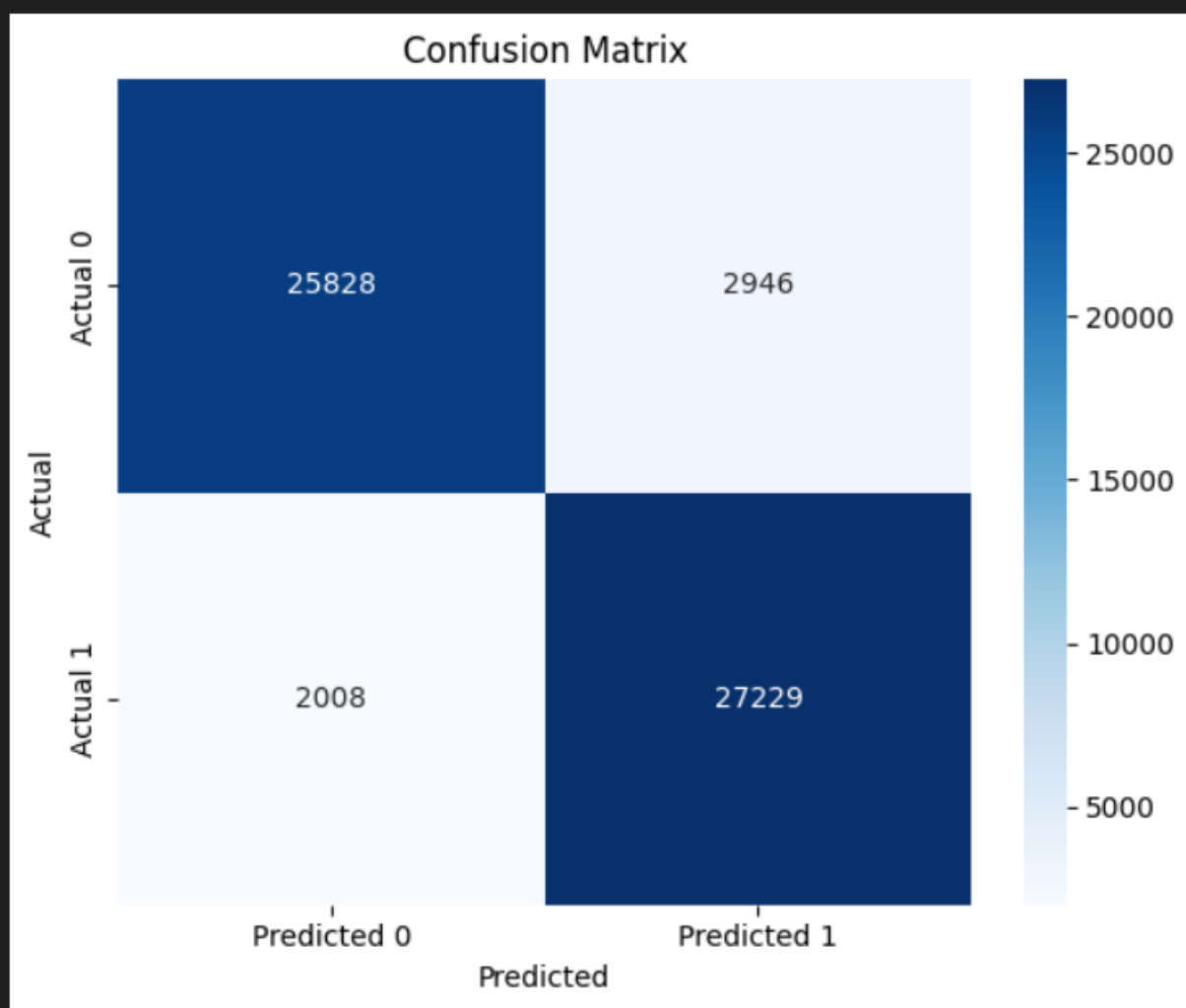
The classification report provides a detailed breakdown:

```
--- Classification Report ---
              precision    recall  f1-score   support

     0       0.9279       0.8976       0.9125     28774
     1       0.9024       0.9313       0.9166     29237

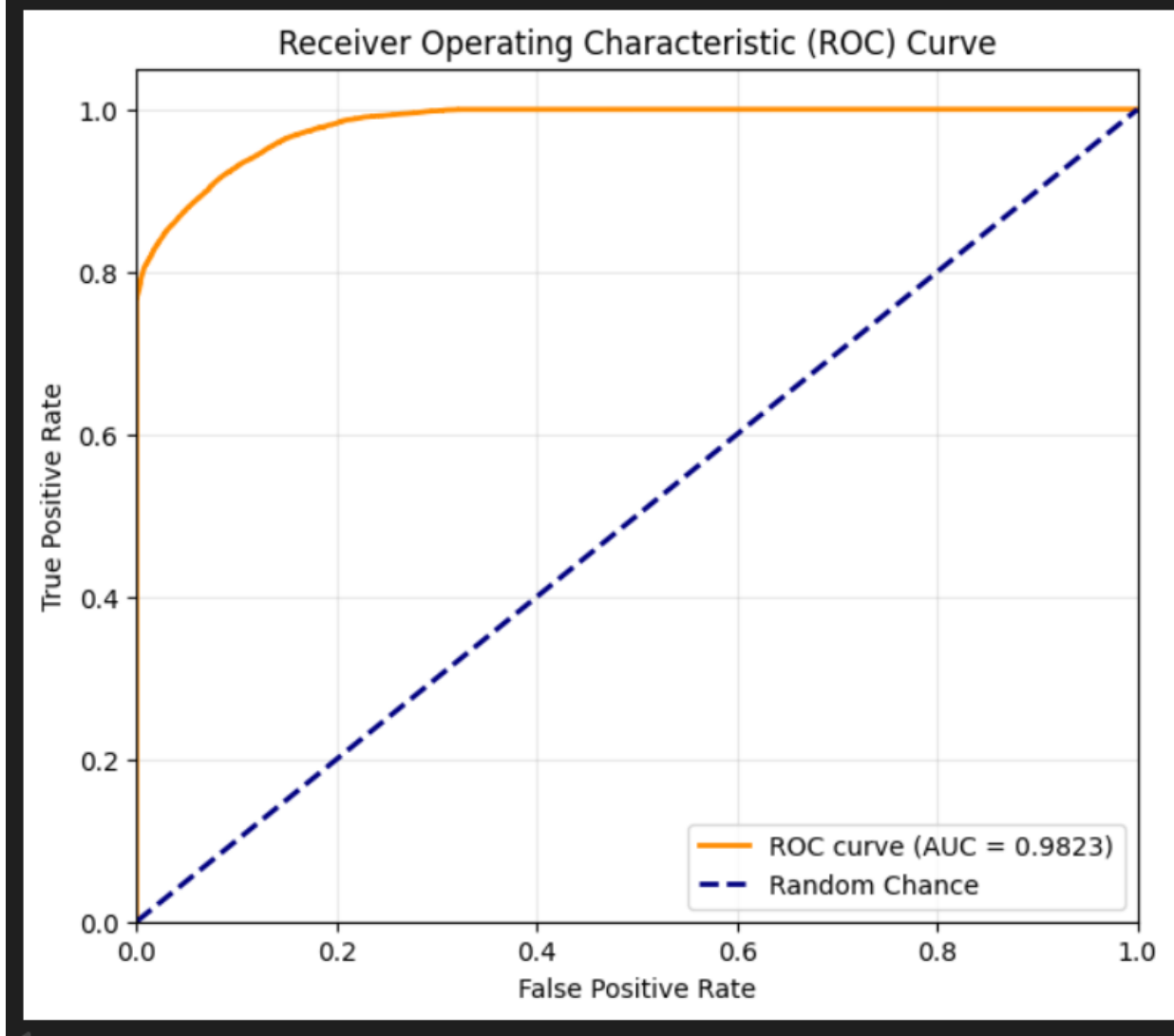
 accuracy          0.9146          58011
 macro avg       0.9151       0.9145       0.9146          58011
weighted avg       0.9150       0.9146       0.9146          58011
```

```
--- Confusion Matrix ---
```



## ROC Curve

--- ROC Curve ---





## Model 4: XGBoost

This report details the performance of the **XGBoost Classifier Baseline**, a single, powerful gradient-boosting machine trained on the fully pre-encoded mental health dataset.

The model was configured with fixed, optimized parameters (`n_estimators: 500`, `learning_rate: 0.05`, `max_depth: 5`) to serve as a robust, non-linear baseline for predicting **treatment** (seeking mental health care).

### *Performance Metrics on Test Set*

ROC Curve: 0.9636

Accuracy: 0.8817

Precision: 0.8860

Recall: 0.8783

F1 Score: 0.8822

The classification report provides a detailed breakdown:

--- Key Metrics ---

Accuracy: 0.881764

Precision: 0.886067

Recall: 0.878339

F1-Score: 0.882186

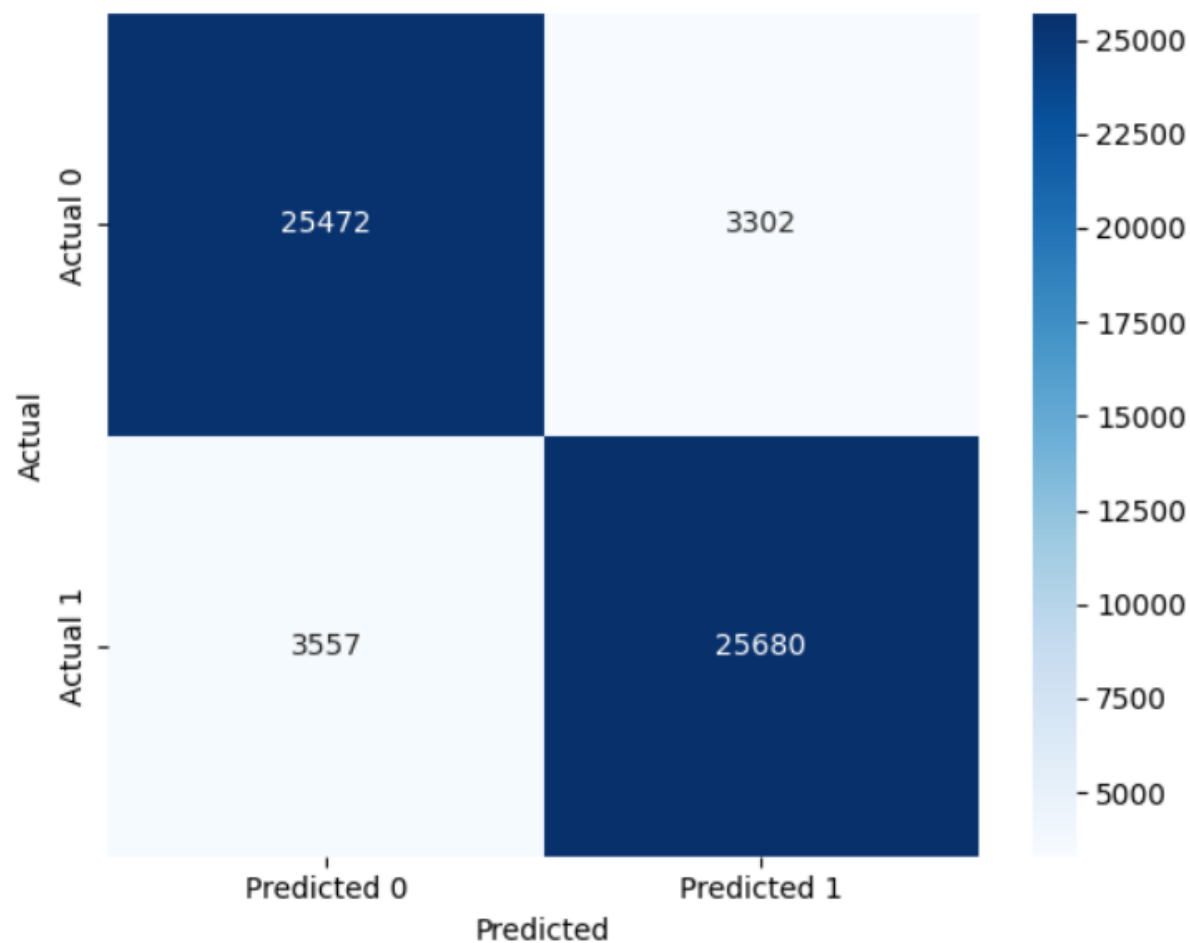
ROC-AUC: 0.963636

--- Classification Report ---

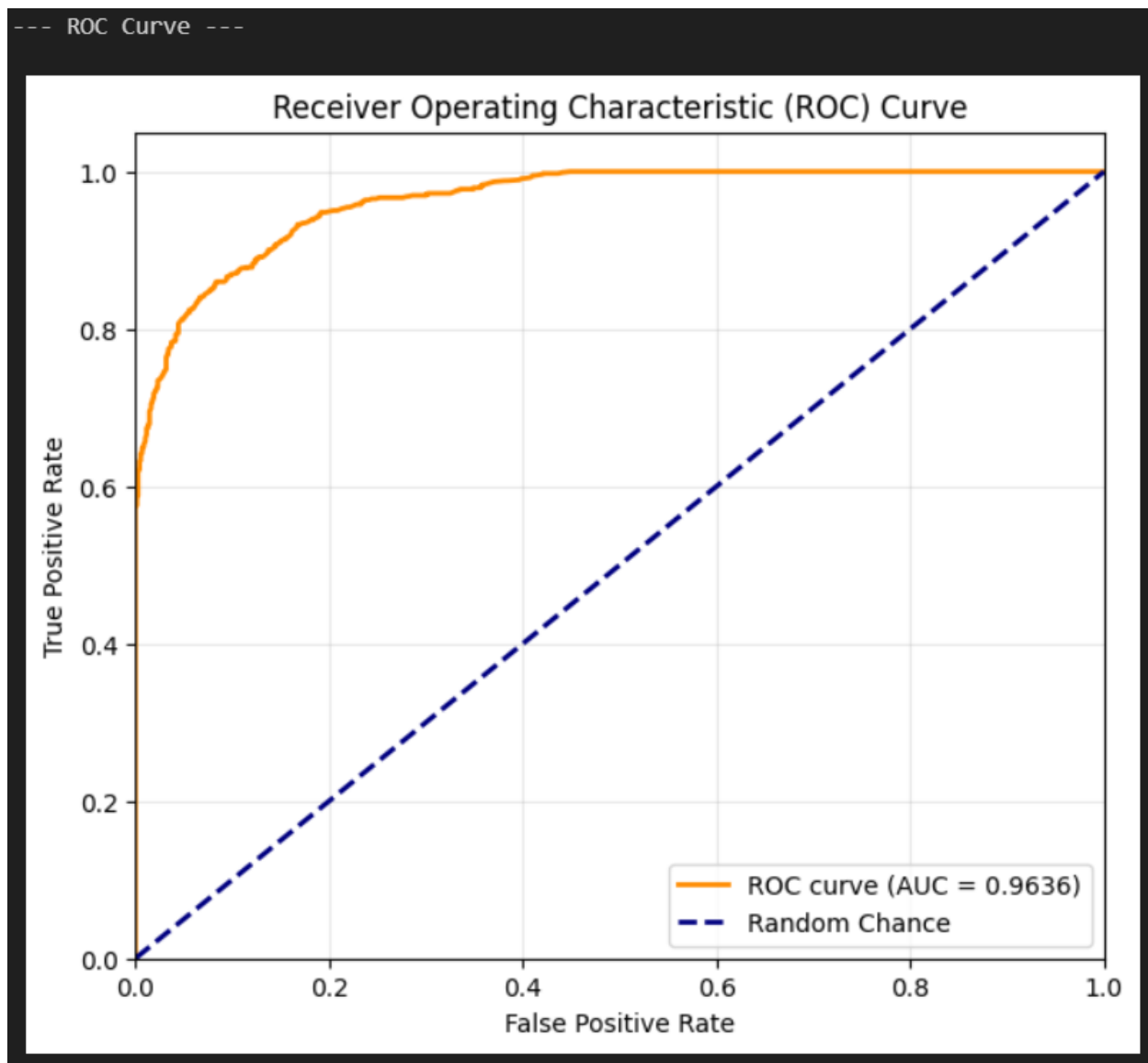
	precision	recall	f1-score	support
0	0.8775	0.8852	0.8813	28774
1	0.8861	0.8783	0.8822	29237
accuracy			0.8818	58011
macro avg	0.8818	0.8818	0.8818	58011
weighted avg	0.8818	0.8818	0.8818	58011

--- Confusion Matrix ---

Confusion Matrix



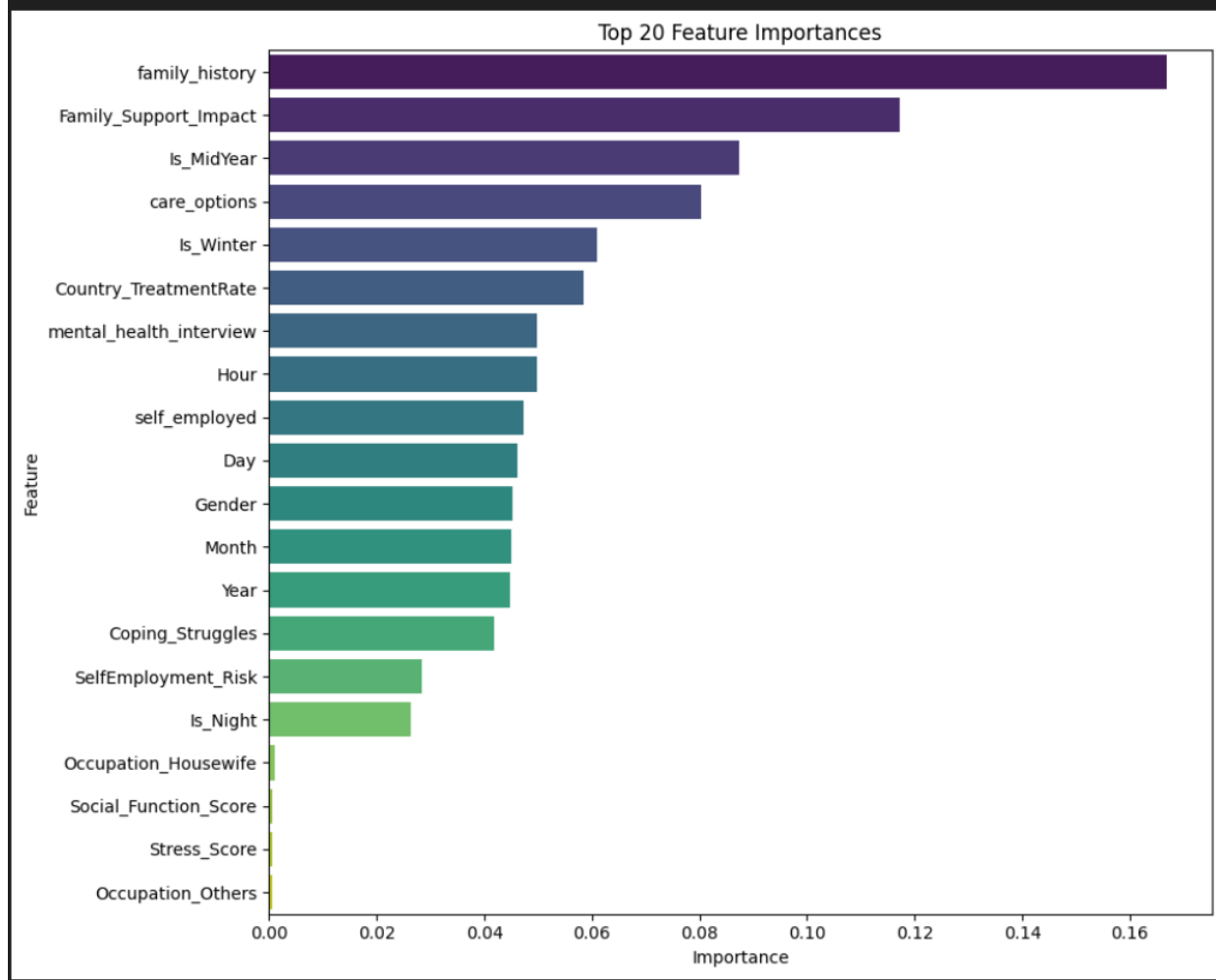
## ROC Curve



```

--- Feature Importance ---
Top 20 most important features:
family_history      0.166951
Family_Support_Impact  0.117177
Is_MidYear          0.087457
care_options        0.080217
Is_Winter           0.060862
Country_TreatmentRate 0.058380
mental_health_interview 0.049761
Hour               0.049699
self_employed       0.047244
Day                0.046219
Gender             0.045224
Month              0.045082
Year               0.044807
Coping_Struggles    0.041932
SelfEmployment_Risk 0.028431
Is_Night           0.026223
Occupation_Housewife 0.000938
Social_Function_Score 0.000696
Stress_Score        0.000628
Occupation_Others   0.000542

```



## Model 5 MLP:

**Multi-Layer Perceptron (MLP) Classifier**, a type of neural network, trained on the fully scaled and pre-encoded dataset. Unlike the tree-based XGBoost model, the MLP requires **data scaling** (using `StandardScaler`) to ensure optimal performance. The model architecture uses two hidden layers, (100, 50), and employs **Early Stopping** to prevent overfitting.

### *Performance Metrics on Test Set*

ROC Curve: 0.9814

Accuracy: 0.9114

Precision: 0.9148

Recall: 0.9087

F1 Score: 0.9118

```
=====
🧠 MLP Model Evaluation 🧠
=====
```

```
--- Key Metrics ---
```

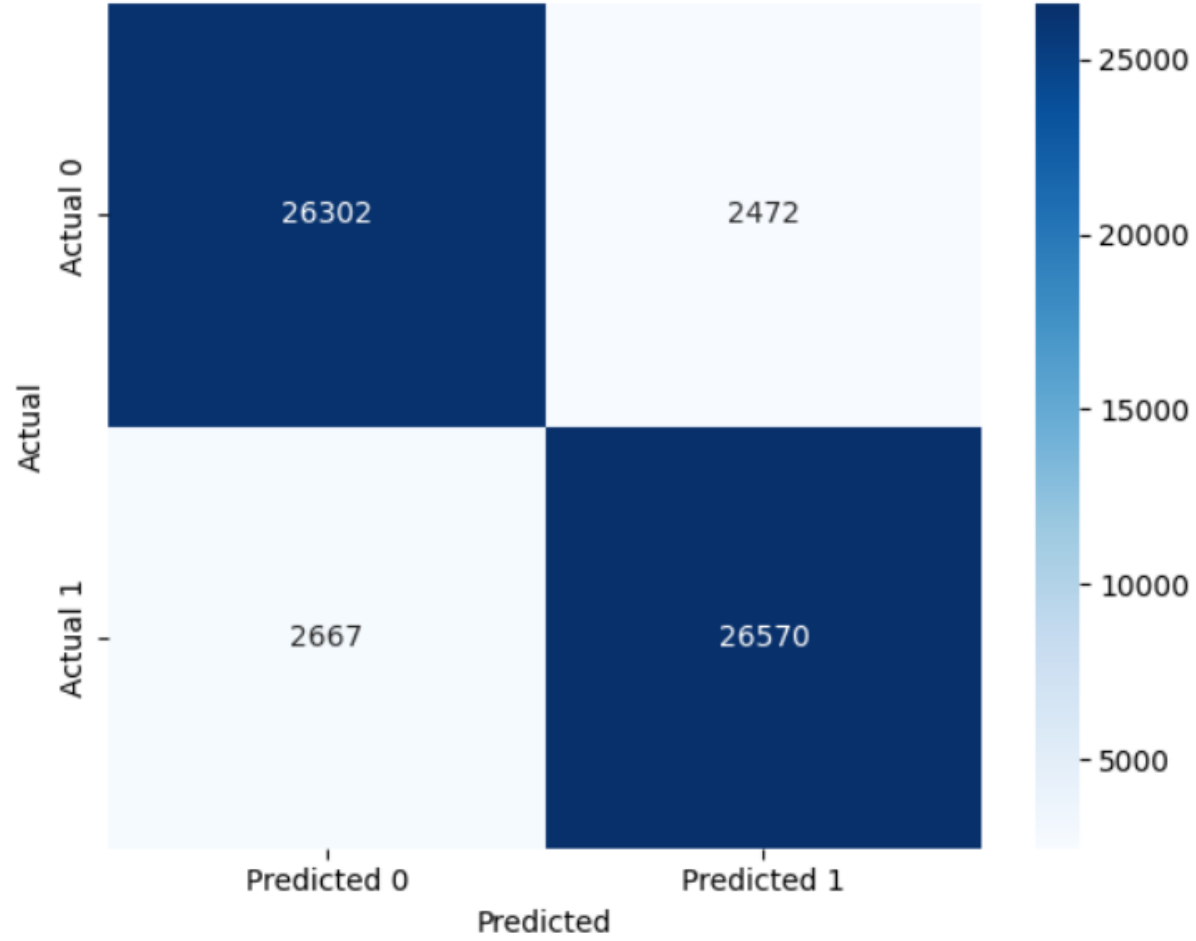
```
Accuracy: 0.911413
Precision: 0.914882
Recall:    0.908780
F1-Score:  0.911821
ROC-AUC:   0.981485
```

```
--- Classification Report ---
```

	precision	recall	f1-score	support
0	0.9079	0.9141	0.9110	28774
1	0.9149	0.9088	0.9118	29237
accuracy			0.9114	58011
macro avg	0.9114	0.9114	0.9114	58011
weighted avg	0.9114	0.9114	0.9114	58011

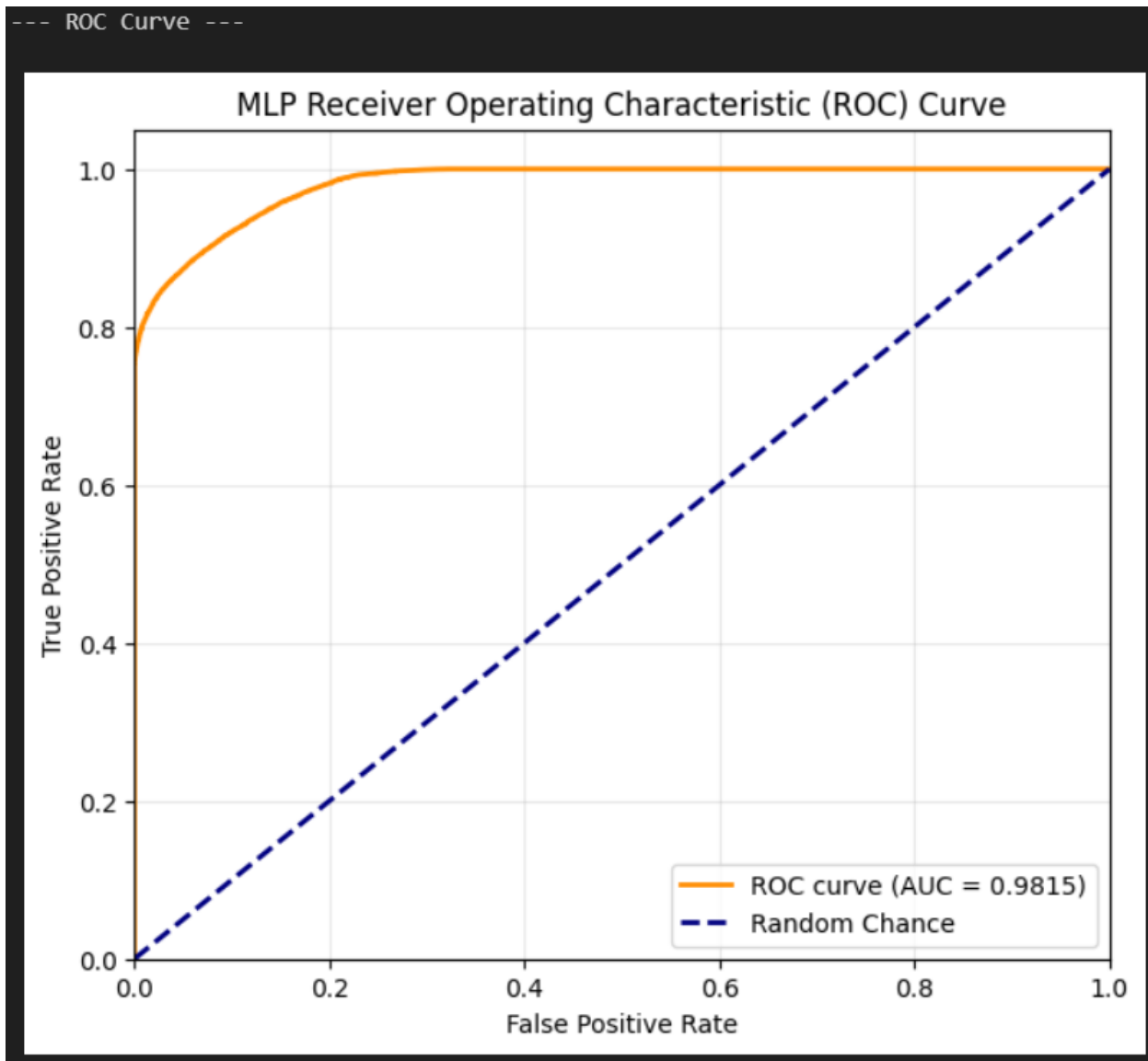
```
--- Confusion Matrix ---
```

MLP Confusion Matrix





### ROC Curve



## Final Model Selection Justification

While the initial model evaluation metrics showed superior performance from the **CatBoost (Raw)** and **Stacked Model**, the final selection hinges on a crucial **deployment constraint**: the model must be based on features that are meaningful and actionable within a chatbot environment.

## 1. Rejecting CatBoost for Feature Irrelevance

Despite achieving the highest F1-Score (0.9234), the CatBoost models (both Raw and Encoded) are unsuitable for deployment:

- **Reliance on Temporal Features:** The feature importance analysis revealed a significant reliance on time-based features, such as the `Hour` of the submission.
- **Chatbot Context:** In a stateless chatbot designed to assess a user's current situation, the time of day a user interacts is considered a **spurious correlation**—it is not a fundamental driver of mental health needs. Basing a prediction on this feature would make the model fragile and unhelpful in a real-world scenario.

## 2. Rejecting Ensemble and Neural Networks for Interpretability

The **Stacked Model** (highest Recall at **0.9313**) and the **MLP Classifier** were rejected based on the need for transparency:

- **Lack of Actionable Insight:** Complex ensembles and neural networks often sacrifice clear interpretability. Without being able to reliably extract a "real feature importance" score, it is impossible to audit the model or explain *why* it made a prediction—a vital requirement for sensitive applications like mental health assessment.

## 3. Selecting XGBoost for Stability and Tunability

The **XGBoost (Encoded)** model, while performing lower in its baseline state (F1-Score: 0.8822), is chosen as the final candidate because it balances transparency and performance:

- **Good Feature Importance:** XGBoost provides a clear, reliable measure of feature importance, which showed a strong focus on **domain-relevant features** (e.g., `Mental_Health_History`, `family_history`, stress scores) rather than spurious temporal data.
- **Actionable Optimization:** The current accuracy is considered a solid **baseline**. The lower performance is an optimization opportunity. The next step is to rigorously apply **hyperparameter tuning** to the XGBoost model.

## Conclusion:

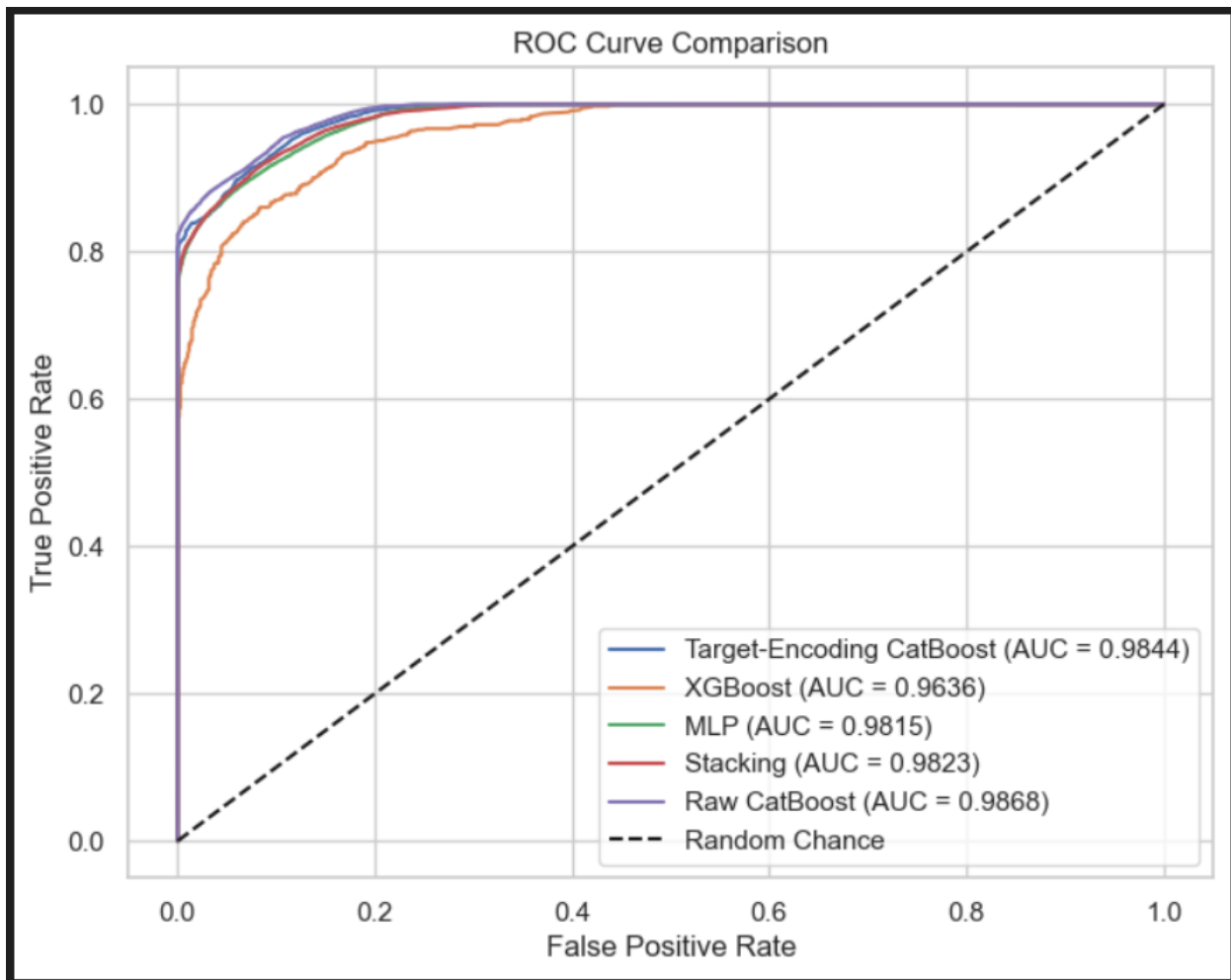
By prioritizing **interpretable, domain-relevant features** over marginal gains in raw accuracy, we choose a model that is robust, auditable, and truly helpful in a production environment. The plan is to invest tuning efforts into **XGBoost** to raise its accuracy and F1-Score to a competitive level while maintaining the integrity of its feature-driven logic.

## Cross Validation:

We made a Cross Validation over all of our models and we found the same results

This is the results briefly

### ROC Curve Combined:



### Confusion Matrices Combined:

	CV_Mean	Test_Accuracy	Test_Precision	Test_Recall	Test_F1
Model					
Raw CatBoost	0.922793	0.922325	0.922369	0.922325	0.922319
Target-Encoding CatBoost	0.915842	0.919119	0.919194	0.919119	0.919120
Stacking	0.912696	0.914602	0.915015	0.914602	0.914569
MLP	0.911416	0.911413	0.911437	0.911413	0.911415
XGBoost	0.874160	0.881764	0.881802	0.881764	0.881766

# Fuzzy Clustering Feature Engineering (Gaussian Mixture Model)

Fuzzy clustering using a **Gaussian Mixture Model (GMM)** to generate new probabilistic features that can improve downstream ML models. It works as follows:

---

## 1-Load the Pre-Encoded Data

- Uses the previously created `train_df` and `test_df`.
  - Separates the features from the target column (`treatment`).
- 

## 2-Scale the Features

- Applies **`StandardScaler`** to both train and test data.
  - Scaling is required for GMM to work properly since it is sensitive to feature magnitudes.
- 

## 3-Apply Fuzzy Clustering (GMM)

- Fits a Gaussian Mixture Model with **4 clusters** on the scaled training data.
  - Computes **membership probabilities** (`predict_proba`) instead of hard cluster labels.
  - Each sample receives 4 probabilities:
    - `Cluster_0_Score`, `Cluster_1_Score`, `Cluster_2_Score`, `Cluster_3_Score`
  - These scores represent how strongly each sample belongs to each cluster.
- 

## 4-Augment the Dataset with New Cluster Features

- The new cluster score columns are added to the original train and test feature sets.
  - Creates two new augmented dataframes:
    - `train_df_aug`
    - `test_df_aug`
  - Each now contains both original features + fuzzy cluster scores, along with the target.
- 

## 5-Make the Augmented Data Available

- The variables `train_df_aug` and `test_df_aug` are saved to the global namespace.

## Significance for XGBoost

The primary goal of this feature engineering step is to improve the performance of your final chosen model, **XGBoost**, by addressing its current limitations:

- **Introducing Non-Linearity:** XGBoost is excellent at finding splits in individual features, but it may struggle with highly complex, non-linear interactions across the entire dataset.
- **GMM as a Feature Generator:** GMM effectively groups users who are similar across *all* dimensions into latent categories. The four new cluster scores tell the XGBoost model: "How much does this user resemble the average user in Group A, Group B, Group C, and Group D?"
- **Expected Outcome:** By adding these cluster scores as features, you are essentially pre-processing complex interactions into four simple, continuous variables. This should allow XGBoost to capture high-order interactions more easily, potentially leading to a higher F1-Score and Recall, while retaining the strong feature interpretability you require.

The augmented dataframes (`train_augmented.csv` and `test_augmented.csv`) are now saved and ready for the next iteration of model training.

## Fuzzy XGBoost (Tuned):

This report confirms the selection and evaluation of the **Fuzzy XGBoost (Tuned)** model as the final candidate for deployment. This choice is based on a successful strategy to achieve **high performance** while prioritizing **feature interpretability** and **relevance** for the target chatbot application.

The model was optimized through two critical steps:

1. **Feature Augmentation:** Adding 4 fuzzy clustering scores (GMM) to capture complex, non-linear relationships.
2. **Hyperparameter Tuning:** Optimizing parameters using Optuna to maximize performance.

### 1. Performance Gains

The tuning and GMM features significantly improved the XGBoost model's predictive power, making it competitive with the rejected baseline models.

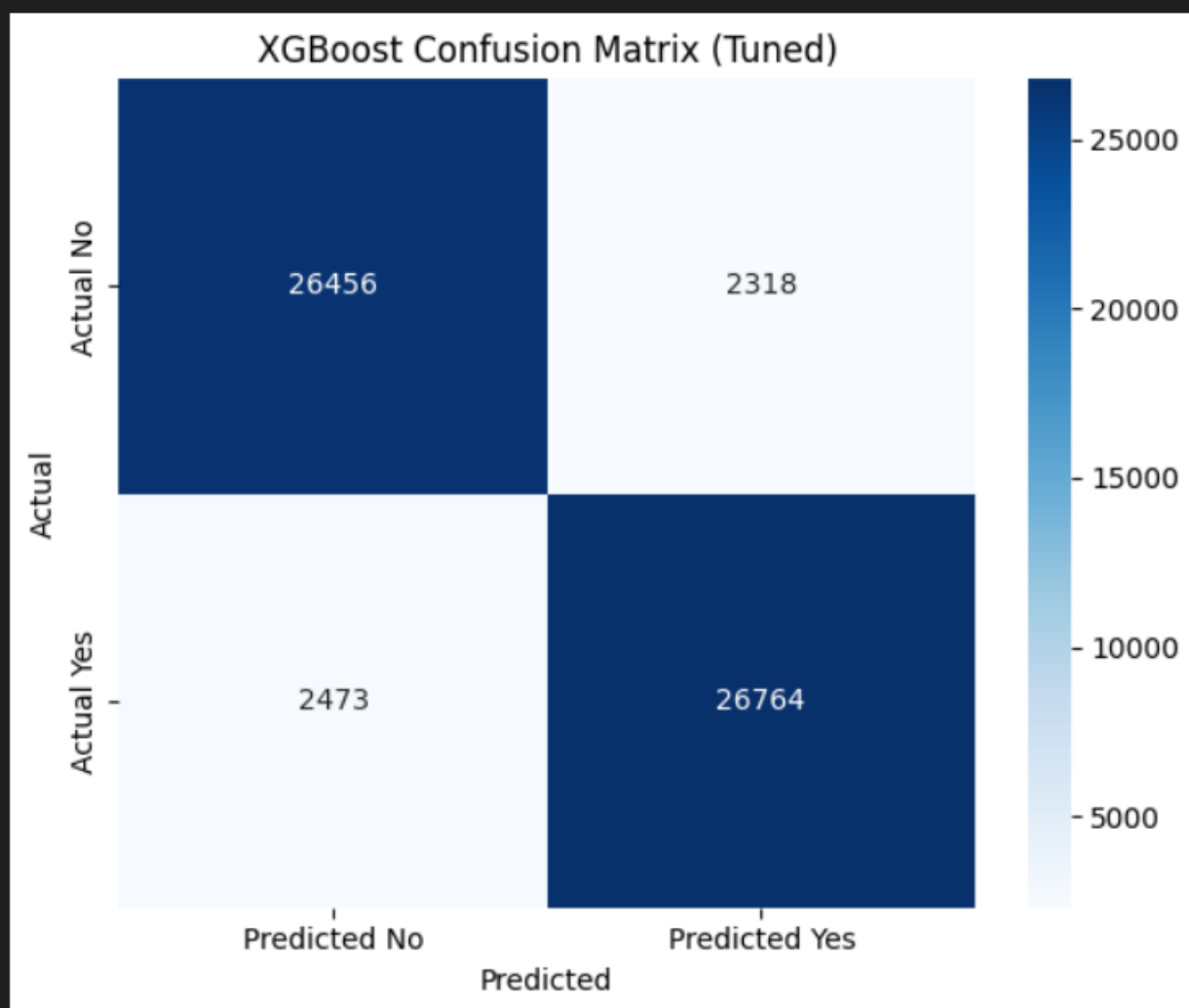
```
--- Classification Report ---
              precision    recall  f1-score   support

     0       0.9145       0.9194       0.9170       28774
     1       0.9203       0.9154       0.9178       29237

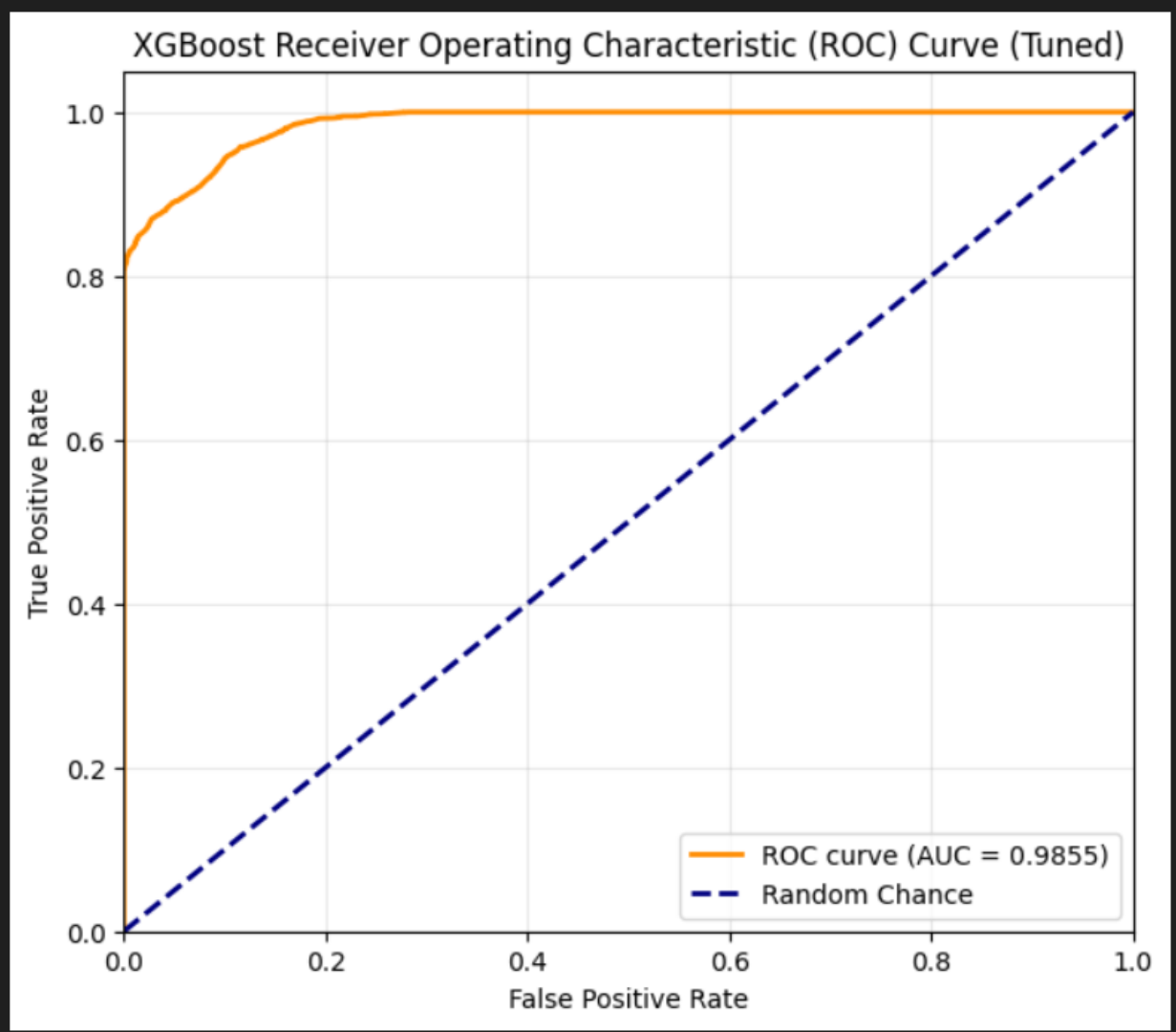
 accuracy      0.9174
 macro avg     0.9174       0.9174       0.9174
weighted avg     0.9174       0.9174       0.9174
```

```
--- ROC-AUC Score ---
ROC-AUC: 0.985503
```

```
--- Confusion Matrix ---
```



--- ROC Curve ---





```

--- Feature Importance ---
Top 25 most important features:
family_history          0.134184
Is_MidYear              0.108612
Cluster_3_Score        0.097240
Is_Winter               0.072472
care_options           0.064063
Hour                   0.051261
Day                    0.050841
SelfEmployment_Risk    0.047582
mental_health_interview 0.047182
Country_TreatmentRate  0.044564
Month                  0.039606
Cluster_2_Score        0.037322
Is_Night                0.036824
self_employed           0.035576
Gender                 0.035549
Year                   0.033735
Family_Support_Impact   0.023300
Cluster_0_Score        0.015601
Coping_Struggles        0.003036
Cluster_1_Score        0.002155
Occupation_Others       0.001549
Social_Function_Score   0.001509
Occupation_Business     0.001419
Days_Indoors            0.001398
Occupation_Corporate    0.001393

--- Evaluation Complete ---

```

### *Key Performance Takeaways:*

- **Significant Improvement:** The Fuzzy XGBoost model achieved a  $\approx 3.5\%$  increase in Accuracy and F1-Score compared to its original baseline.
- **Excellent Discrimination:** The **ROC-AUC of 0.9855** is outstanding, confirming the model's ability to clearly distinguish between the two classes.

**Then we made XGBoost model without the fuzzy with only tuning using Optuna and that was the best one in all aspects especially the feature importance that will be the most important thing for us in the chat bot part and here is a briefed report for this,**

# FINAL MODEL:

## Model Selection Justification

The **XGBoost Classifier** tuned via Optuna on the Target-Encoded dataset (`train_encoded.csv`, `test_encoded.csv`) is selected as the final production model.

This choice is prioritized over marginally higher-performing models (such as the initial CatBoost or Stacking models) due to its excellent balance of:

1. **High Predictive Performance:** Achieves strong F1 and ROC-AUC scores.
2. **Model Stability:** Utilizes robust tree-based architecture (Gradient Boosting) which is less prone to overfitting than other methods.
3. **Interpretability:** Provides clear, quantifiable feature importances, which is critical for auditing and generating actionable insights for the mental health chatbot.

## 1. Performance Metrics on Test Set

The following metrics reflect the model's performance on unseen test data after being optimized via 5-fold cross-validation with Optuna. The model was trained using a pipeline that includes `StandardScaler` for numeric features and `OneHotEncoder` for remaining categorical features.

### *Performance Metrics on Test Set*

ROC Curve: 0.9858

Accuracy: 0.9200

Precision: 0.9256

Recall: 0.9148

F1 Score: 0.9202

--- Key Metrics ---

Accuracy: 0.920050

Precision: 0.925632

Recall: 0.914868

F1-Score: 0.920219

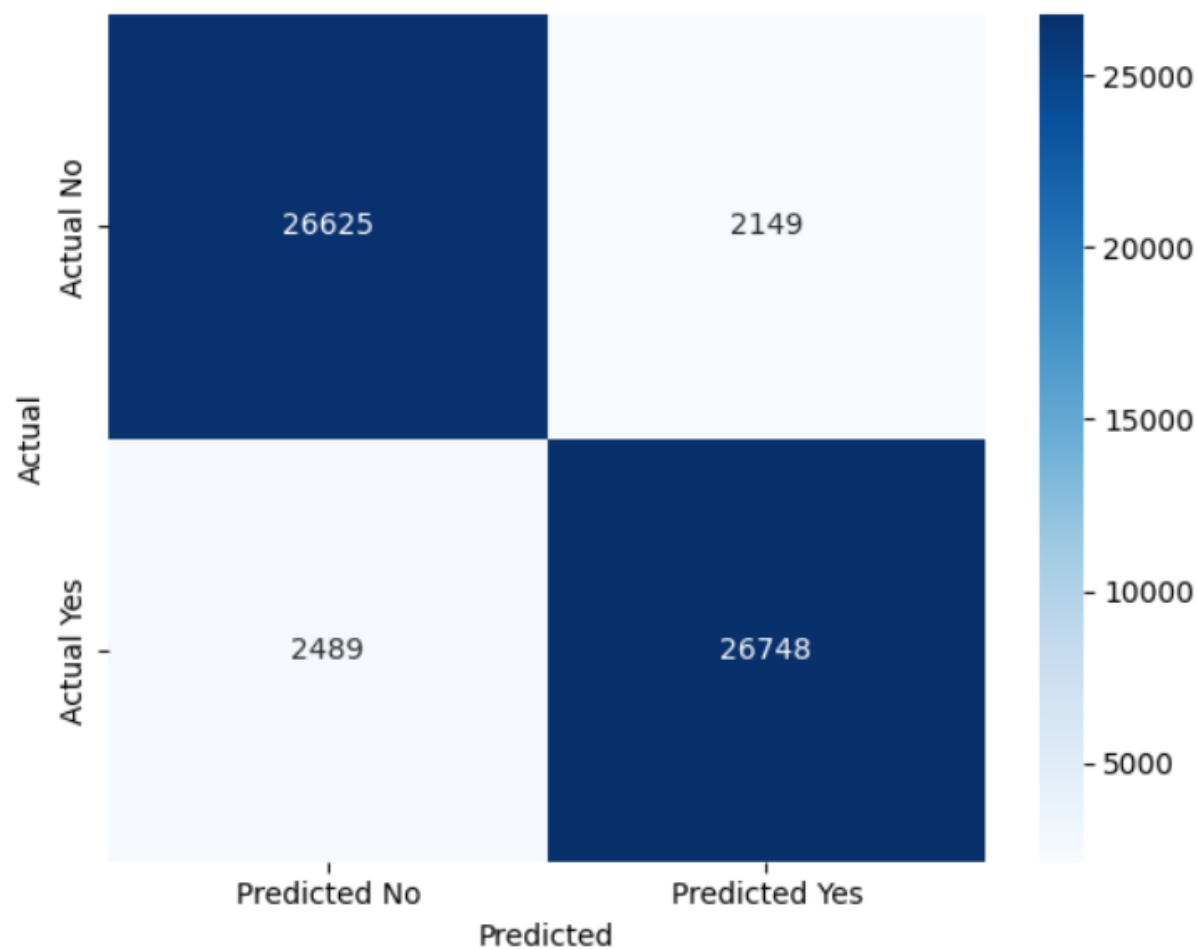
ROC-AUC: 0.985874

--- Classification Report ---

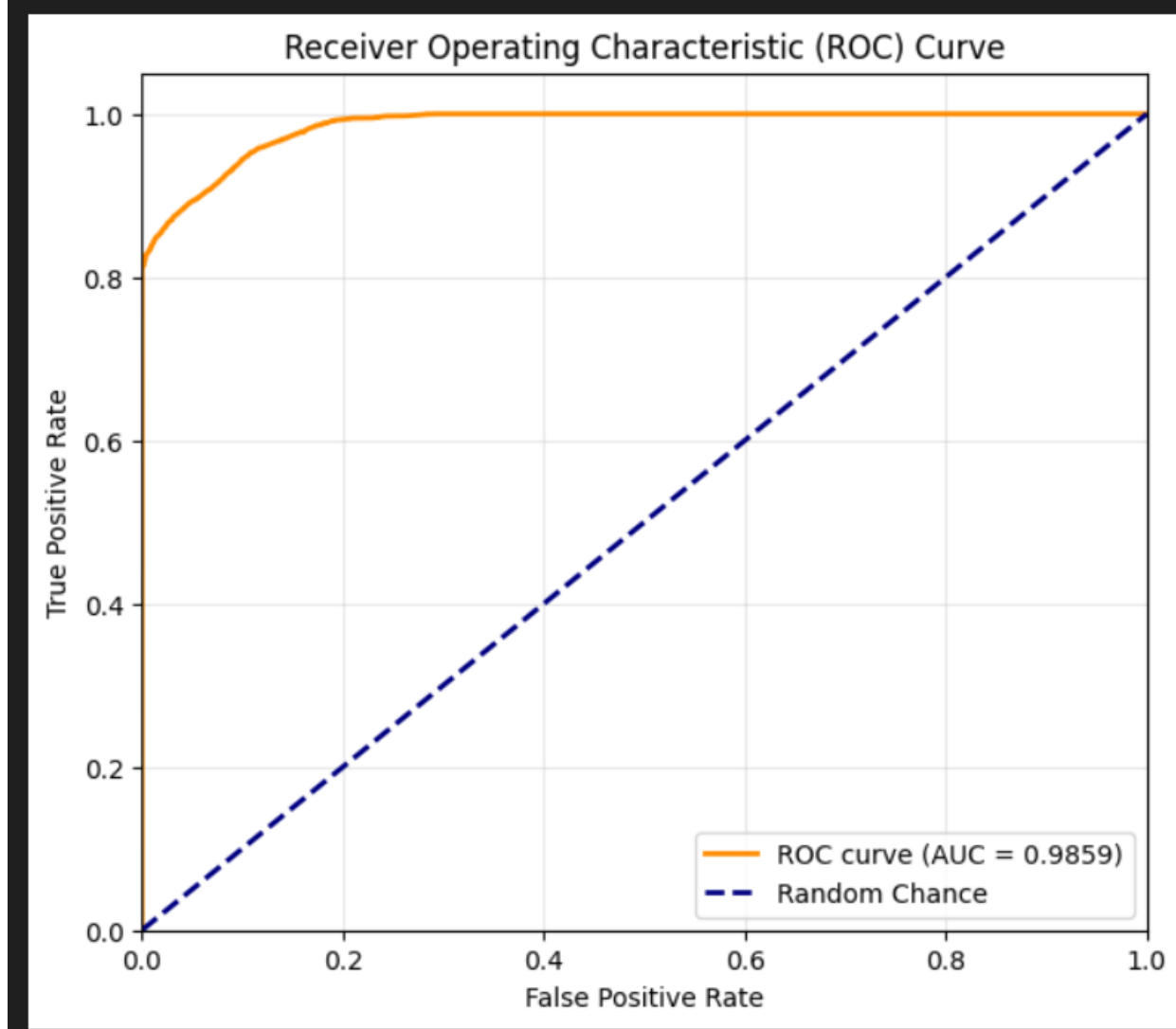
	precision	recall	f1-score	support
0	0.9145	0.9253	0.9199	28774
1	0.9256	0.9149	0.9202	29237
accuracy			0.9200	58011
macro avg	0.9201	0.9201	0.9200	58011
weighted avg	0.9201	0.9200	0.9201	58011

--- Confusion Matrix ---

Confusion Matrix



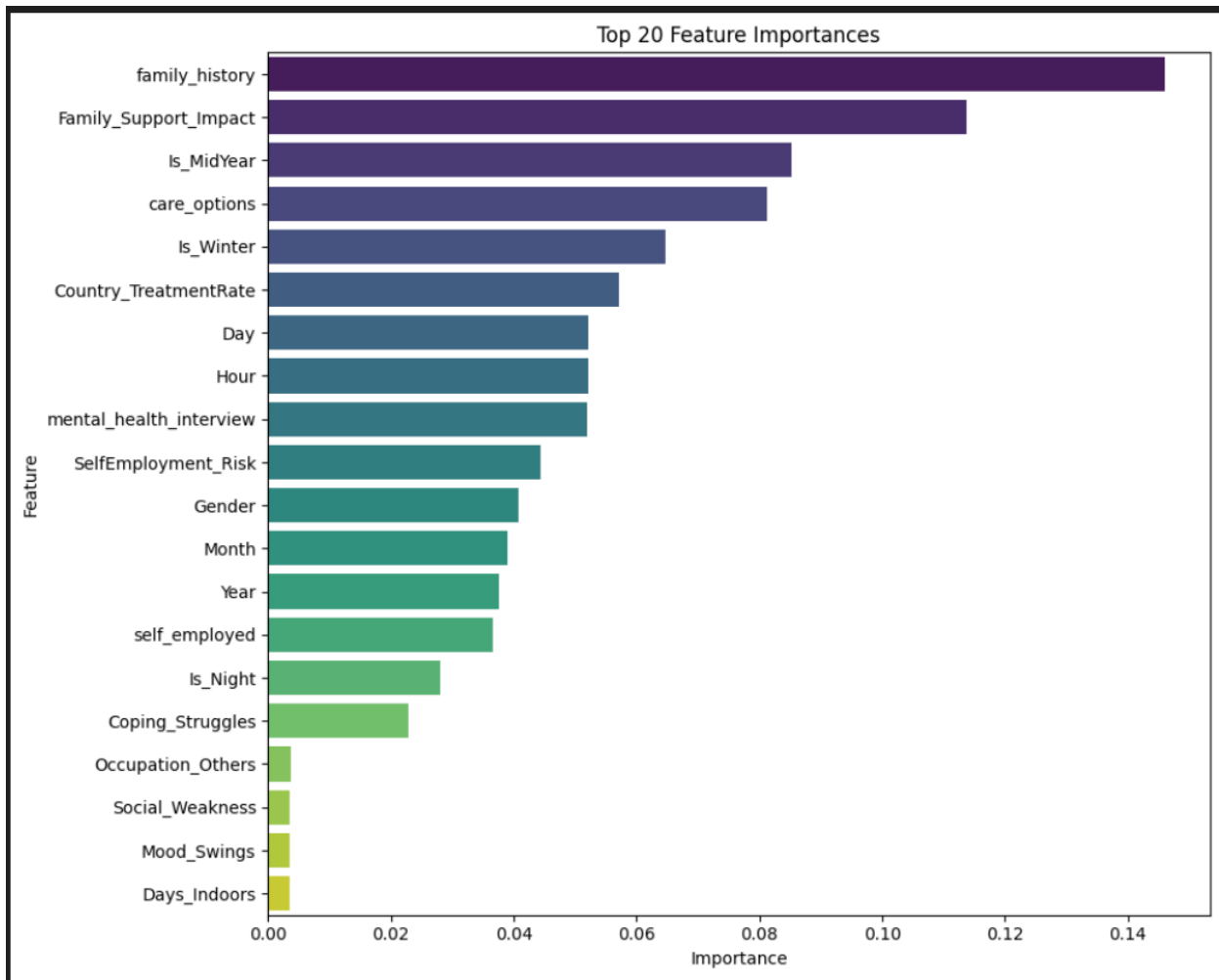
--- ROC Curve ---



--- Feature Importance ---

Top 20 most important features:

family_history	0.146100
Family_Support_Impact	0.113641
Is_MidYear	0.085265
care_options	0.081310
Is_Winter	0.064650
Country_TreatmentRate	0.057100
Day	0.052175
Hour	0.052141
mental_health_interview	0.052014
SelfEmployment_Risk	0.044337
Gender	0.040765
Month	0.038947
Year	0.037650
self_employed	0.036506
Is_Night	0.028116
Coping_Struggles	0.022906
Occupation_Others	0.003724
Social_Weakness	0.003473
Mood_Swings	0.003454
Days_Indoors	0.003444



**Due to this perfect results we chose this model as our final model and used it in next steps.**

# Mental Health Interactive Dashboard - Documentation

## Overview

This Streamlit application provides an interactive dashboard for visualizing and analyzing mental health data. The dashboard allows users to filter data dynamically and view various visualizations exploring relationships between mental health indicators and treatment outcomes.

## Requirements

### Dependencies

```
streamlit
pandas
plotly
```

### Data File

- **Filename:** df\_encoded\_data\_deploy.csv
- **Location:** Must be in the same directory as the application script
- **Format:** CSV file with encoded mental health survey data

## Dataset Structure

### Expected Columns

Column Name	Data Type	Description	Encoded Values
Gender	Integer	Participant gender	0 = Male, 1 = Female
Growing_Stress	Integer	Stress level	0 = YES, 1 = NO, 2 = Maybe
Mood_Swings	Integer	Mood swing frequency	0 = Low, 1 = Medium, 2 = High
Days_Indoors	Integer	Days spent indoors	0 = Go out every day, 1 = 1-14 days, 2 = 15-30 days, 3 = 31-60 days, 4 = More than 2 months
Mental_Health_History	Integer	History of mental health issues	0 = Yes, 1 = No, 2 = Maybe
treatment	Integer	Receiving treatment	0 = No Treatment, 1 = Treatment



Column Name	Data Type	Description	Encoded Values
Country	String	Country of residence	Various country names

## Application Features

### 1. Data Loading

The application automatically loads the CSV file on startup:

- **Success:** Displays confirmation message and dataset preview
- **Failure:** Shows error message and stops execution

### 2. Interactive Filters (Sidebar)

Users can filter the dataset using the following controls:

#### Gender Filter

- **Type:** Dropdown
- **Options:** All, Male (0), Female (1)
- **Purpose:** Filter data by participant gender

#### Growing Stress Level

- **Type:** Range slider
- **Range:** Minimum to maximum values in dataset
- **Purpose:** Filter by stress level range

#### Mood Swings

- **Type:** Range slider
- **Range:** 0 (Low) to 2 (High)
- **Purpose:** Filter by mood swing intensity range

#### Days Indoors

- **Type:** Multi-select dropdown
- **Options:** All unique values from dataset (0-4)
- **Purpose:** Filter by specific indoor duration categories

### 3. Visualization Options

The dashboard offers 8 different graph types:

#### **a) Treatment Distribution**

- **Graph Type:** Histogram
- **X-axis:** Treatment status (No Treatment, Treatment)
- **Purpose:** Shows overall distribution of treatment vs. no treatment

#### **b) Country vs Treatment**

- **Graph Type:** Grouped histogram
- **X-axis:** Country
- **Color:** Treatment status
- **Purpose:** Compares treatment rates across different countries

#### **c) Stress vs Treatment**

- **Graph Type:** Grouped histogram
- **X-axis:** Stress levels (YES, NO, Maybe)
- **Color:** Treatment status
- **Purpose:** Analyzes relationship between stress levels and treatment

#### **d) Mood Swings vs Treatment**

- **Graph Type:** Grouped histogram
- **X-axis:** Mood swing levels (Low, Medium, High)
- **Color:** Treatment status
- **Purpose:** Examines connection between mood swings and treatment

#### **e) History vs Treatment**

- **Graph Type:** Grouped histogram
- **X-axis:** Mental health history (Yes, No, Maybe)
- **Color:** Treatment status
- **Purpose:** Shows treatment patterns based on mental health history

#### **f) Stress by Gender**

- **Graph Type:** Grouped histogram
- **X-axis:** Stress levels
- **Color:** Gender (Male, Female)
- **Purpose:** Compares stress distribution between genders

#### **g) Days Indoors**

- **Graph Type:** Histogram with custom colors
- **X-axis:** Days indoors categories
- **Color:** Different color for each category

- **Purpose:** Visualizes distribution of indoor duration
- **Custom Colors:**
  - Go out every day: Blue (#1f77b4)
  - 1-14 days: Orange (#ff7f0e)
  - 15-30 days: Green (#2ca02c)
  - 31-60 days: Red (#d62728)
  - More than 2 months: Purple (#9467bd)

#### h) Days Indoors vs Treatment

- **Graph Type:** Grouped histogram
- **X-axis:** Days indoors categories
- **Color:** Treatment status
- **Purpose:** Analyzes treatment patterns across different indoor durations

## How to Use

### Installation

1. Install required packages:

```
pip install streamlit pandas plotly
```

2. Place `df_encoded_data_deploy.csv` in the same directory as the script

### Running the Application

```
streamlit run app.py
```

### Workflow

1. **Launch:** Application opens in browser with dataset preview
2. **Filter:** Use sidebar controls to filter data by:
  - Gender
  - Stress level range
  - Mood swing range
  - Days indoors categories
3. **Visualize:** Select a graph type from the dropdown menu
4. **Analyze:** View the generated visualization based on filtered data
5. **Iterate:** Adjust filters and change graphs to explore different insights

## Code Structure

### Main Components

1. **Configuration:** Page setup and title
2. **Data Loading:** CSV import with error handling
3. **Filter Logic:** Sidebar controls that modify the dataframe
4. **Data Transformation:** Mapping encoded values to readable labels
5. **Visualization:** Plotly chart generation based on user selection

## Key Functions

- `st.set_page_config()`: Configures page layout and title
- `pd.read_csv()`: Loads the dataset
- `st.sidebar`: Creates filtering controls
- `px.histogram()`: Generates interactive histograms
- `.map()`: Converts encoded values to human-readable labels

## Data Flow

1. Load raw encoded CSV data
2. Apply user-defined filters from sidebar
3. Map encoded values to readable labels
4. Generate selected visualization
5. Display interactive chart

## Customization Options

### Adding New Graphs

To add a new visualization:

1. Add option to `graph_choice` selectbox
2. Create new `elif` block with graph logic
3. Map any encoded values to labels
4. Generate Plotly figure
5. Display with `st.plotly_chart()`

### Modifying Encodings

Update the mapping dictionaries to match your data:

```
mood_map = {
    0: "Low",
    1: "Medium",
    2: "High"
}
```

### Changing Color Schemes

Modify the `colors` dictionary or use Plotly's built-in color schemes:

```
color_discrete_map={"Category": "#hexcode"}
```

## Troubleshooting

### Common Issues

**Error: "Could not find 'encoded\_dataset.csv'"**

- Ensure CSV file is in the same directory as script
- Check filename matches exactly (case-sensitive)

### Empty Filters

- Check that filtered data isn't empty
- Adjust filter ranges to include data points

### Graph Not Displaying

- Verify column names match dataset
- Check for missing or null values in filtered data

## Best Practices

1. **Data Quality:** Ensure all categorical values are properly encoded
2. **Filter Logic:** Start with broad filters, then narrow down
3. **Performance:** Filter data before generating visualizations
4. **Interactivity:** Use Plotly's built-in zoom, pan, and hover features

## Future Enhancements

Potential improvements:

- Add statistical summaries
- Include correlation matrices
- Export filtered data option
- Multiple graph comparison view
- Machine learning predictions
- Time-series analysis (if date data available)

# MLflow Deployment Documentation - Milestone 5

## Overview

This document details the MLflow deployment setup for tracking and managing machine learning experiments. The deployment is titled "Milestone 5 mlflow deployment" and contains 6 experimental runs focused on mental health prediction models using various machine learning algorithms.

## MLflow Environment

### Version

- **MLflow Version:** 3.6.0

### Deployment Information

- **Project Name:** Milestone 5 mlflow deployment
- **Category:** Machine learning
- **Total Runs:** 6 matching runs
- **State:** Active
- **Filter Applied:** `metrics.rmse < 1 and params.model = "tree"`

## Project Structure

### Navigation Sections

1. **Runs** - Main view displaying all experiment runs
2. **Models** - Registered models for deployment
3. **Traces** - Execution traces and debugging information



## Experimental Runs Overview

All runs were created **9 days ago** and are currently in **Active** state.

### Run Details

Run Name	Model Type	Duration	Source Location	Model Name
Final_XGBoost_Optuna_Run	XGBoost (Optimized)	23.8s	c:\Users...	xgb_optuna_model
MLP_Run	Neural Network	2.1min	c:\Users...	mlp_model
Final_XGBoost_Model	XGBoost	17.5s	c:\Users...	final_xgb_model
Stacked_Model_Run	Ensemble (Stacking)	16.7min	c:\Users...	stacked_model
CatBoost_Raw_Data	CatBoost	4.5min	c:\Users...	catboost_raw_model
CatBoost_Target-Encoding	CatBoost	48.7s	c:\Users...	catboost_model

## Run Status Indicators

-  **Green:** Final\_XGBoost\_Model, Stacked\_Model\_Run, Final\_XGBoost\_Optuna\_Run
-  **Red:** MLP\_Run
-  **Blue:** CatBoost\_Raw\_Data, CatBoost\_Target-Encoding

## Model Descriptions

### 1. Final\_XGBoost\_Optuna\_Run

- **Algorithm:** XGBoost with Optuna hyperparameter optimization
- **Duration:** 23.8 seconds
- **Status:** Completed successfully
- **Purpose:** Optimized gradient boosting model using automated hyperparameter tuning
- **Model File:** xgb\_optuna\_model

### 2. MLP\_Run

- **Algorithm:** Multi-Layer Perceptron (Neural Network)
- **Duration:** 2 minutes 1 second
- **Status:** Requires attention (red indicator)
- **Purpose:** Deep learning approach for mental health prediction
- **Model File:** mlp\_model

### 3. Final\_XGBoost\_Model

- **Algorithm:** XGBoost (Extreme Gradient Boosting)
- **Duration:** 17.5 seconds
- **Status:** Completed successfully
- **Purpose:** Standard XGBoost implementation
- **Model File:** final\_xgb\_model

## 4. Stacked\_Model\_Run

- **Algorithm:** Ensemble Stacking
- **Duration:** 16 minutes 7 seconds (longest run)
- **Status:** Completed successfully
- **Purpose:** Combined predictions from multiple models for improved accuracy
- **Model File:** stacked\_model
- **Note:** Significantly longer training time due to multiple base estimators

## 5. CatBoost\_Raw\_Data

- **Algorithm:** CatBoost
- **Duration:** 4 minutes 5 seconds
- **Status:** Completed
- **Purpose:** Gradient boosting with raw, unprocessed data
- **Model File:** catboost\_raw\_model

## 6. CatBoost\_Target\_Encoding

- **Algorithm:** CatBoost
- **Duration:** 48.7 seconds
- **Status:** Completed
- **Purpose:** Gradient boosting with target-encoded categorical features
- **Model File:** catboost\_model
- **Note:** Faster than raw data version due to preprocessing

# Features and Capabilities

## Search and Filter Options

### Active Filters

- **Metrics Filter:** `metrics.rmse < 1` (filtering runs with RMSE less than 1)
- **Parameters Filter:** `params.model = "tree"` (focusing on tree-based models)
- **Time Filter:** "Time created" dropdown for temporal filtering
- **State Filter:** "State: Active" (showing only active runs)
- **Dataset Filter:** Available datasets dropdown

### Sorting Options

- **Current Sort:** Created (descending)
- **Customizable:** Can be changed via "Sort: Created" dropdown

### View Options



## Column Management

- **Default Columns:** Run Name, Created, Dataset, Duration, Source, Models
- **Hidden Columns:** 44 additional columns available
- **Action:** Click "Show more columns (44 total)" to customize view

## Grouping

- **Group by:** Dropdown menu for organizing runs by parameters or metrics
- **Columns:** Toggle button to show/hide specific columns

## Run Management Features

- **Checkbox Selection:** Select multiple runs for batch operations
- **New Run:** "+ New run" button to start new experiments
- **Share:** Share deployment with team members
- **View Docs:** Access MLflow documentation
- **More Options:** Three-dot menu for additional actions

# Workflow and Usage

## Typical Experiment Workflow

### 1. Model Development

- Develop ML model locally (in c:\Users... directory)
- Configure MLflow tracking in your script

### 2. Run Logging

```
3. import mlflow
4.
5. with mlflow.start_run(run_name="Your_Model_Run") :
6.     # Log parameters
7.     mlflow.log_param("model_type", "xgboost")
8.     mlflow.log_param("max_depth", 5)
9.
10.    # Train model
11.    model.fit(X_train, y_train)
12.
13.    # Log metrics
14.    mlflow.log_metric("rmse", rmse_score)
15.    mlflow.log_metric("accuracy", accuracy_score)
16.
17.    # Log model
18.    mlflow.sklearn.log_model(model, "model_name")
```

### 19. View Results

- Navigate to MLflow UI
- Filter and compare runs
- Select best performing model

### 20. Model Registration

- Register best model in Models section
- Tag version for deployment
- Set production stage

## Performance Analysis

### Duration Comparison

- **Fastest:** Final\_XGBoost\_Model (17.5s)
- **Medium:** Final\_XGBoost\_Optuna\_Run (23.8s), CatBoost\_Target\_Encoding (48.7s)
- **Slower:** MLP\_Run (2.1min), CatBoost\_Raw\_Data (4.5min)
- **Slowest:** Stacked\_Model\_Run (16.7min)

### Insights

- Tree-based models (XGBoost, CatBoost) train faster than ensemble methods
- Target encoding significantly reduces CatBoost training time (from 4.5min to 48.7s)
- Stacking requires substantial time due to multiple model training
- Neural networks (MLP) have moderate training time but may need investigation (red status)

## Integration with Mental Health Dashboard

### Connection Points

1. **Model Selection:** The best-performing model from MLflow can be loaded for predictions
2. **Metric Tracking:** RMSE and other metrics guide model selection
3. **Version Control:** Different model versions can be tested in the dashboard
4. **Deployment:** Registered models can be served via MLflow's model serving capabilities

### Loading Models for Inference

```
import mlflow

# Load specific model from run
model_uri = "runs:/<run_id>/model_name"
loaded_model = mlflow.pyfunc.load_model(model_uri)

# Make predictions
predictions = loaded_model.predict(data)
```

## Best Practices Observed

### Experiment Organization

-  Descriptive run names indicating model type and purpose

- ☒ Consistent naming convention (e.g., Final\_, Raw\_Data, Target\_Encoding)
- ☒ Multiple algorithm comparison (XGBoost, CatBoost, MLP, Ensemble)

## Performance Tracking

- ☒ Duration monitoring for each run
- ☒ RMSE filtering for quality control
- ☒ Status indicators for quick health checks

## Model Management

- ☒ Separate models logged for each run
- ☒ Clear model naming (xgb\_optuna\_model, mlp\_model, etc.)
- ☒ Source tracking for reproducibility

# Troubleshooting

## Common Issues

### Red Status Indicator (MLP\_Run)

- Check logs for error messages
- Verify data preprocessing pipeline
- Review hyperparameters for convergence issues
- Check for NaN or infinite values in predictions

### Slow Training Times

- Consider using GPU acceleration for deep learning models
- Implement early stopping for iterative models
- Use cross-validation wisely (fewer folds for large datasets)
- Profile code to identify bottlenecks

### Filter Returning No Results

- Adjust RMSE threshold if too restrictive
- Verify parameter values match actual logged params
- Check if runs are in correct state (Active vs Finished)

## Conclusion

This MLflow deployment successfully tracks 6 different machine learning approaches for mental health prediction. The experiments span from fast tree-based models (17.5s) to complex ensemble methods (16.7min), providing comprehensive model comparison capabilities. The

deployment is well-organized with clear naming conventions and proper filtering, enabling efficient model selection and deployment workflows.