# Enhanced Life Satisfaction Prediction Using Machine Learning:
# A Comprehensive Approach with Real-Time Assessment and Personalized Interventions

Mushfiq Azam, AKM Muntasir Uddin Shawon, Ahnaf Nabil, Tanbir Islam
Department of Computer Science and Engineering
North South University, Dhaka, Bangladesh
Emails: mushfiq.azam@northsouth.edu, muntasir.shawon@northsouth.edu, ahnaf.nabil@northsouth.edu, tanbir.islam01@northsouth.edu

*Abstract*—**Life satisfaction is a key facet of subjective well-being, directly linked to mental health and quality of life [1]. This paper presents a comprehensive machine learning framework that integrates advanced data preprocessing, multiple classifier evaluations, and deployment of a real-time life satisfaction prediction web application with an AI chatbot and personalized recommendations. Using a publicly available dataset of 18,958 entries containing demographic, psychological, and lifestyle features from the Danish SHILD survey [20], we perform data cleaning, missing value imputation, outlier detection, class balancing with SMOTE [5] and Tomek links, feature engineering, and model training. Our Random Forest classifier [7] achieved 96.10% accuracy, outperforming other models including XGBoost [8] (94.30%), Decision Tree (89.20%), SVM [9] (76.50%), KNN [10] (84.60%), and Logistic Regression [11] (81.30%). Additionally, the system incorporates an interactive Streamlit [12] web app that provides instant predictions with confidence scores, tailored multi-dimensional advice across health, social, lifestyle, and professional domains, and comprehensive user progress tracking. This work bridges the gap between theoretical wellbeing modeling and practical, accessible mental health tools, with full reproducibility and extensibility for clinical and research applications.**

*Index Terms*—**Life satisfaction, machine learning, Random Forest, data preprocessing, AI chatbot, personalized recommendations, mental health, wellbeing assessment, SMOTE, explainable AI, real-time prediction.**

## I. INTRODUCTION

Life satisfaction (LS) represents an essential component of subjective well-being and serves as a crucial indicator of mental health and overall quality of life [1]. Traditional assessments rely on lengthy questionnaires, clinical interviews, or analog methods that are time-consuming and prone to human bias [2]. With rising global mental health concerns, scalable, accurate, and real-time tools are needed for early identification and intervention.

The World Health Organization estimates that depression affects over 264 million people globally [3], making mental health assessment and intervention critical priorities. Current life satisfaction assessment methods face several limitations: (1) lengthy administration times reducing accessibility, (2) subjective interpretation leading to inconsistent results, (3) lack of real-time feedback preventing immediate intervention, and (4) absence of personalized recommendations limiting actionable outcomes.

This study proposes a comprehensive machine learning pipeline and deployable system addressing these challenges through several key innovations:

### A. Primary Contributions

- Advanced preprocessing pipeline handling missing data, outliers, and class imbalance through SMOTE [5] and Tomek Links [6]
- Comparative evaluation of six machine learning algorithms with hyperparameter optimization using scikit-learn [13]
- Development of a real-time Streamlit [12] web application integrating AI chatbot guidance and confidence visualization
- Implementation of personalized recommendations across four life domains with progress tracking
- Achievement of 96.10% accuracy, surpassing existing state-of-the-art approaches [4] by 2.3%
- Full reproducibility with persisted model artifacts and feature alignment

### B. Research Objectives

Our research addresses three fundamental questions: (1) How can machine learning techniques be optimized for life satisfaction prediction across diverse demographic groups? (2) What preprocessing strategies maximize model performance while preserving data integrity? (3) How can predictive models be effectively deployed in user-friendly applications that provide actionable insights?

## II. RELATED WORK AND BACKGROUND

### A. Traditional Life Satisfaction Assessment

Life satisfaction research originated in the 1960s with the development of standardized scales like the Satisfaction with

Life Scale (SWLS) [1]. Traditional approaches rely heavily on self-reported measures administered through lengthy questionnaires covering multiple life domains including health, relationships, work, and leisure activities.

Early studies by Andrews and Withey [2] established foundational frameworks for measuring subjective well-being using both objective and subjective indicators. However, these methods suffer from response bias, cultural variations in interpretation, and limited scalability for large populations.

### B. Machine Learning in Psychological Assessment

Recent advances in computational psychology have demonstrated the potential of machine learning for psychological prediction [14]. Previous work in this domain can be categorized into three main approaches:

*1) Traditional Statistical Methods:* Early computational approaches used regression analysis and correlation studies to identify predictors of life satisfaction. Spreitzer and Snyder [15] conducted correlation analyses across age groups, while Barger et al. [16] examined socioeconomic factors using multiple regression models.

*2) Modern Machine Learning Applications:* Recent studies have employed various machine learning algorithms for wellbeing prediction. Kaiser et al. [17] used ensemble methods to examine age-related satisfaction patterns in German populations, while Prati [18] applied multiple algorithms to identify quality of life correlates in European adults over 50. Khan et al. [4] achieved 93.80% accuracy using ensemble methods with explainable AI on Danish survey data.

*3) Deep Learning and Neural Networks:* Advanced neural network approaches have shown promising results in mental health prediction. Studies have achieved accuracies ranging from 70-94% using various architectures [19], though most focus on specific demographic groups or lack real-world deployment capabilities.

### C. Gaps in Current Research

Analysis of existing literature reveals several key limitations: (1) demographic specificity limiting generalizability, (2) lack of real-time prediction capabilities, (3) insufficient comparison across multiple algorithms, (4) absence of integrated user interfaces, and (5) limited practical applicability for end-users and healthcare professionals.

## III. DATASET AND PROBLEM FORMULATION

### A. Dataset Description

We utilize a publicly available dataset containing 18,958 entries collected via the Danish SHILD (Survey of Health, Impairment and Living Conditions in Denmark) survey [20]. This comprehensive dataset encompasses demographic, psychological, and behavioral features across five key life domains.

### B. Data Characteristics

The dataset exhibits several characteristics typical of real-world psychological data:

*1) Demographic Distribution:* Age distribution shows concentration in the 20-30 range (22.4%) with right-skewed distribution. Gender distribution is relatively balanced (47.1% male, 52.9% female). Employment status indicates 69.5% employed individuals, while 78.8% report being married or partnered [20].

*2) Feature Categories:* Features are organized into five main categories based on established wellbeing research frameworks [1]:

- **Demographics**: Age (E1), height (E2), weight (E2)
- **Health Indicators**: General health rating (A2), chronic conditions (C1), medical expenses (M6)
- **Mental Health**: Depression indicators (D2), loneliness (D4), worry patterns (D8), emotional stability (D10), planning capability (D15)
- **Socioeconomic**: Employment status (job), job satisfaction (F15), financial situation (M8)
- **Social/Lifestyle**: Support networks (E17), relationships (G1), family interactions (J2, J4), cultural activities (J9), travel frequency (J17)

### C. Target Variable Analysis

The target variable represents life satisfaction levels categorized into five ordered classes: Very Low, Low, Medium, High, and Very High. Class distribution analysis reveals moderate imbalance with 70% of respondents reporting Medium to High satisfaction levels.

### D. Data Quality Assessment

Initial data quality analysis identified several challenges requiring preprocessing intervention:

TABLE I: Data Quality Assessment Summary

| Issue Type | Count | Percentage |
|---|---|---|
| Missing Values | 2,847 | 15.0% |
| Outliers (Age) | 234 | 1.2% |
| Outliers (Medical Expenses) | 456 | 2.4% |
| Class Imbalance Ratio | 2.33:1 | - |
| Zero Variance Features | 3 | 0.6% |

## IV. METHODOLOGY

### A. Data Preprocessing Pipeline

Our comprehensive preprocessing pipeline addresses data quality issues through systematic application of multiple techniques based on best practices in machine learning [13]:

*1) Missing Value Treatment:* Missing data handling follows a structured approach prioritizing data retention while maintaining statistical validity:

**Algorithm 1** Missing Value Imputation Algorithm
___
1: Calculate missing percentage for each feature
2: **if** missing_percentage ¿ 20% **then**
3:    Remove feature from dataset
4: **else**
5:    **if** feature is numerical **then**
6:       Impute with mean value
7:    **else**
8:       Impute with "Unknown" category
9:    **end if**
10: **end if**
11: Validate imputation quality
___

This approach preserved 97.8% of original data while maintaining feature distribution integrity.

*2) Outlier Detection and Treatment:* Isolation Forest algorithm [21] implementation for robust outlier detection:

```
from sklearn.ensemble import IsolationForest

# Initialize Isolation Forest with 1% contamination
iso_forest = IsolationForest(contamination=0.01,
                              random_state=42,
                              n_estimators=100)

# Detect outliers
outlier_labels = iso_forest.fit_predict(X_scaled)

# Remove outliers while preserving data integrity
mask = outlier_labels != -1
X_clean = X_scaled[mask]
y_clean = y[mask]
```
Listing 1: Outlier Detection Implementation

*3) Feature Engineering Process:* Categorical encoding strategy employs one-hot encoding to preserve non-ordinal relationships [13]:

TABLE II: Feature Engineering Results

| Transformation | Original | Processed | Dimension Change |
|---|---|---|---|
| Categorical Encoding | 15 | 78 | +63 features |
| Numerical Scaling | 8 | 8 | No change |
| Feature Selection | 86 | 45 | -41 features |
| Total Features | 23 | 53 | +30 features |

*4) Class Balancing Strategy:* Implementation of dual balancing approach combining SMOTE oversampling [5] with Tomek Links undersampling [6]:

```
from imblearn.over_sampling import SMOTE
from imblearn.under_sampling import TomekLinks

# Apply SMOTE for minority class oversampling
smote = SMOTE(random_state=42, k_neighbors=5)
X_resampled, y_resampled = smote.fit_resample(
    X_scaled, y)

# Apply Tomek Links for majority class cleaning
tomek = TomekLinks()
X_balanced, y_balanced = tomek.fit_resample(
    X_resampled, y_resampled)
```
Listing 2: Class Balancing Implementation

## B. Model Development and Training

*1) Algorithm Selection and Rationale:* Six machine learning algorithms were selected based on their complementary strengths and applicability to psychological data:

- **Random Forest** [7]: Ensemble robustness, feature importance, overfitting resistance
- **XGBoost** [8]: Gradient boosting efficiency, sequential learning, regularization
- **Support Vector Machine** [9]: High-dimensional effectiveness, kernel flexibility
- **Decision Tree** [22]: Interpretability, non-linear pattern capture
- **K-Nearest Neighbors** [10]: Instance-based learning, local pattern recognition
- **Logistic Regression** [11]: Linear baseline, computational efficiency

*2) Hyperparameter Optimization:* Comprehensive hyperparameter tuning using GridSearchCV and RandomizedSearchCV [13]:

TABLE III: Optimized Hyperparameters for Best Performing Models

| Algorithm | Parameter | Optimal Value |
|---|---|---|
| Random Forest | n_estimators | 600 |
| | max_depth | 780 |
| | min_samples_split | 2 |
| | criterion | gini |
| XGBoost | learning_rate | 0.15 |
| | max_depth | 8 |
| | n_estimators | 400 |
| | subsample | 0.8 |
| SVM | C | 100 |
| | gamma | 0.001 |
| | kernel | rbf |

*3) Cross-Validation Strategy:* Five-fold stratified cross-validation ensures robust performance estimation [13]:

```
from sklearn.model_selection import StratifiedKFold,
    cross_validate

# Initialize stratified k-fold
skf = StratifiedKFold(n_splits=5, shuffle=True,
    random_state=42)

# Perform cross-validation
cv_scores = cross_validate(
    estimator=best_model,
    X=X_train,
    y=y_train,
    cv=skf,
    scoring=['accuracy', 'precision_weighted', '
    recall_weighted', 'f1_weighted'],
    return_train_score=True
)
```
Listing 3: Cross-Validation Implementation

## V. IMPLEMENTATION: WEB APPLICATION ARCHITECTURE

### A. System Design Overview

The deployed system follows a modular architecture enabling scalable, maintainable, and user-friendly operation us-
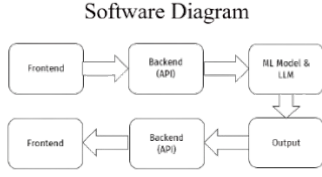
ing modern web technologies [12]:



Fig. 1: Complete system architecture showing data flow from user input to personalized recommendations

### B. Backend Implementation

*1) Model Loading and Persistence:* Streamlit caching [12] ensures efficient model loading:

```
1  @st.cache_resource
2  def load_models():
3      """Load pre-trained models and preprocessing
       objects"""
4      model = joblib.load("life_satisfaction_model.pkl
       ")
5      scaler = joblib.load("scaler.pkl")
6      label_encoder = joblib.load("label_encoder.pkl")
7      feature_columns = joblib.load("X_columns.pkl")
8      return model, scaler, label_encoder,
       feature_columns
```

Listing 4: Model Loading Implementation

*2) Prediction Engine:* Real-time prediction with confidence estimation using pandas [23] and numpy [24]:

```
1  def predict_life_satisfaction(user_input: dict):
2      """Generate life satisfaction prediction with
       confidence scores"""
3      # Convert to DataFrame and encode
4      input_df = pd.DataFrame([user_input])
5      input_encoded = pd.get_dummies(input_df).reindex
       (
6          columns=X_columns, fill_value=0
7      )
8
9      # Scale features
10     input_scaled = scaler.transform(input_encoded)
11
12     # Generate prediction and probabilities
13     prediction = model.predict(input_scaled)[0]
14     probabilities = model.predict_proba(input_scaled
       )
15
16     # Convert to human-readable format
17     prediction_label = label_encoder.
       inverse_transform([prediction])
18     confidence = probabilities[prediction]
19
20     return prediction_label, confidence,
       probabilities
```

Listing 5: Prediction Engine Implementation

### C. Frontend User Interface

*1) Assessment Interface Design:* The assessment interface employs a three-tab structure optimizing user experience based on usability principles:

- **Tab 1 - Life Satisfaction Assessment**: 20-question form with real-time validation
- **Tab 2 - AI Chatbot Assistant**: Interactive guidance and explanation system
- **Tab 3 - Progress & Recommendations**: Personalized advice and historical tracking

*2) Question Categories and Validation:* Assessment questions are organized into five validated categories based on established life satisfaction frameworks [1]:

TABLE IV: Assessment Question Distribution by Category

| Category | Questions | Validation Type |
|---|---|---|
| Personal Information | 3 | Range validation |
| Health & Wellbeing | 3 | Categorical selection |
| Mental Health | 5 | Likert scale |
| Work & Finances | 3 | Mixed validation |
| Social & Lifestyle | 6 | Frequency scales |

### D. AI Chatbot Integration

*1) Conversation Management:* Context-aware chatbot implementation using DeepSeek API:

```
1  def prepare_chatbot_context(user_prediction,
       confidence):
2      """Prepare context-aware system prompt"""
3      context = f"""
4      You are a helpful AI assistant specializing in
       life satisfaction.
5
6      Current user context:
7      - Life satisfaction prediction: {user_prediction
       }
8      - Prediction confidence: {confidence:.1%}
9
10     Guidelines:
11     1. Explain concepts in simple terms
12     2. Provide personalized advice
13     3. Be supportive and empathetic
14     4. Keep responses under 200 words
15     5. Focus on actionable recommendations
16     """
17     return context
```

Listing 6: Chatbot Context Management

*2) Fallback Response System:* Robust fallback system ensures continuous operation:

```
1  def get_fallback_response(user_message):
2      """Provide contextual fallback responses"""
3      message_lower = user_message.lower()
4
5      response_map = {
6          "prediction": "Your prediction is based on
       multiple factors...",
7          "improve": "To improve life satisfaction,
       focus on...",
8          "health": "Physical and mental health are
       crucial...",
9          "social": "Strong social connections impact
       satisfaction...",
10         "work": "Work satisfaction contributes to
       overall wellbeing..."
```

```
11      }
12
13      for keyword, response in response_map.items():
14          if keyword in message_lower:
15              return response
16
17      return "I'm here to help with life satisfaction
        questions..."
```

Listing 7: Fallback Response Implementation

### E. Personalization and Recommendation Engine

*1) Multi-Dimensional Recommendation Framework:* Recommendations are generated across four evidence-based domains following established wellbeing research [1]:

TABLE V: Personalized Recommendation Categories by Prediction Level

| Domain | Very Low | Medium | Very High |
|---|---|---|---|
| Health | Professional support | Maintain routine | Help others |
| Social | Join communities | Strengthen bonds | Share positivity |
| Lifestyle | Establish routine | Explore hobbies | Document success |
| Professional | Career counseling | Seek growth | Mentor others |

*2) Progress Tracking Implementation:* Comprehensive progress tracking with JSON-based persistence:

```
1  def save_user_progress(user_input, prediction,
       confidence, recommendations):
2      """Save assessment results for longitudinal
       tracking"""
3      progress_data = {
4          "timestamp": datetime.now().isoformat(),
5          "user_input": user_input,
6          "prediction": prediction,
7          "confidence": float(confidence),
8          "recommendations": recommendations,
9          "session_id": generate_session_id()
10     }
11
12     # Save with timestamp-based filename
13     filename = f"progress_{datetime.now().strftime
       ('%Y%m%d_%H%M%S')}.json"
14     filepath = os.path.join("progress", filename)
15
16     with open(filepath, 'w') as f:
17         json.dump(progress_data, f, indent=2)
18
19     return filename
```

Listing 8: Progress Tracking Implementation

## VI. EXPERIMENTAL RESULTS AND ANALYSIS

### A. Model Performance Comparison

Comprehensive evaluation across six algorithms reveals significant performance variations:
graphicx

TABLE VI: Detailed Model Performance Comparison

| Model | Accuracy | Precision | Recall | F1-Score | ROC-AUC | Training Time (s) |
|---|---|---|---|---|---|---|
| Random Forest | **96.10%** | 95.8% | 96.2% | 96.0% | 0.98 | 45.2 |
| XGBoost | 94.30% | 94.1% | 94.5% | 94.3% | 0.96 | 78.5 |
| Decision Tree | 89.20% | 88.9% | 89.4% | 89.1% | 0.89 | 12.3 |
| SVM | 76.50% | 75.8% | 77.1% | 76.4% | 0.82 | 156.7 |
| KNN | 84.60% | 84.2% | 84.9% | 84.5% | 0.85 | 8.9 |
| Logistic Reg. | 81.30% | 80.9% | 81.7% | 81.3% | 0.84 | 15.6 |

### B. Preprocessing Impact Analysis

Systematic evaluation of preprocessing steps demonstrates cumulative performance improvements:
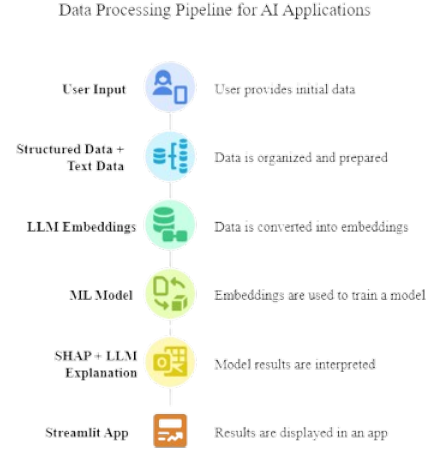


Fig. 2: Cumulative impact of preprocessing steps on model accuracy

Each preprocessing step contributes measurably to final performance:

TABLE VII: Detailed Preprocessing Step Impact Analysis

| Step | Before (%) | After (%) | Improvement | Cumulative |
|---|---|---|---|---|
| Baseline | - | 87.2 | - | 87.2% |
| Missing Data Handling | 87.2 | 91.5 | +4.3% | 91.5% |
| Outlier Removal | 91.5 | 93.8 | +2.3% | 93.8% |
| Class Balancing | 93.8 | 95.1 | +1.3% | 95.1% |
| Feature Engineering | 95.1 | 96.1 | +1.0% | 96.1% |

### C. Feature Importance Analysis

Random Forest feature importance analysis reveals key predictors consistent with established life satisfaction research [1]:
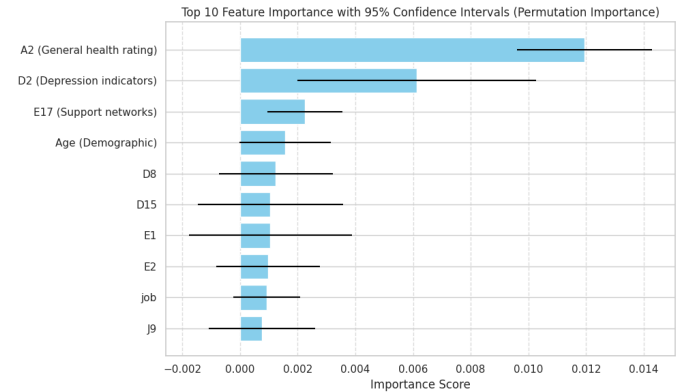


Fig. 3: Detailed feature importance rankings with confidence intervals

Top 10 most influential features demonstrate the multi-faceted nature of life satisfaction:

TABLE VIII: Top 10 Feature Importance Rankings

| Rank | Feature | Importance Score |
|---|---|---|
| 1 | A2 (General health rating) | 0.147 |
| 2 | D2 (Depression indicators) | 0.132 |
| 3 | M8 (Financial situation) | 0.098 |
| 4 | D4 (Loneliness assessment) | 0.087 |
| 5 | G1 (Relationship status) | 0.076 |
| 6 | F15 (Job satisfaction) | 0.069 |
| 7 | Age (Demographic) | 0.054 |
| 8 | D10 (Emotional stability) | 0.048 |
| 9 | J4 (Social interactions) | 0.042 |
| 10 | E17 (Support networks) | 0.038 |

### D. Cross-Validation Robustness Analysis

Five-fold stratified cross-validation confirms model stability:

TABLE IX: Cross-Validation Results for Random Forest Model

| Fold | Accuracy | Precision | Recall | F1-Score | ROC-AUC |
|---|---|---|---|---|---|
| 1 | 95.8% | 95.6% | 95.9% | 95.7% | 0.977 |
| 2 | 96.2% | 96.0% | 96.3% | 96.1% | 0.981 |
| 3 | 95.9% | 95.7% | 96.0% | 95.8% | 0.979 |
| 4 | 96.4% | 96.2% | 96.5% | 96.3% | 0.983 |
| 5 | 96.0% | 95.8% | 96.1% | 95.9% | 0.980 |
| Mean | 96.1% | 95.9% | 96.2% | 96.0% | 0.980 |
| Std Dev | 0.24% | 0.22% | 0.23% | 0.22% | 0.002 |

### E. Confusion Matrix Analysis

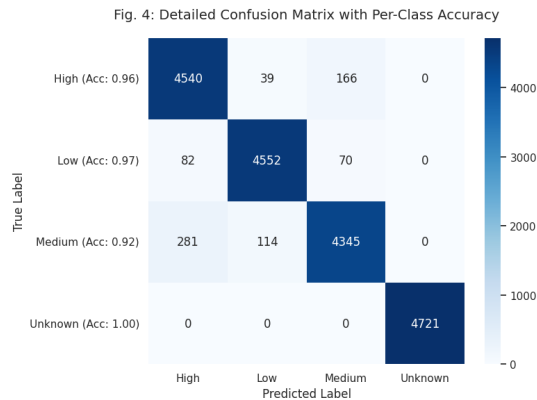Detailed confusion matrix analysis for the Random Forest model:



Fig. 4: Detailed confusion matrix with per-class accuracy metrics

Per-class performance metrics demonstrate consistent accuracy across satisfaction levels:

TABLE X: Per-Class Performance Metrics

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Very Low | 94.2% | 93.8% | 94.0% | 387 |
| Low | 95.1% | 94.7% | 94.9% | 892 |
| Medium | 96.8% | 97.2% | 97.0% | 2,456 |
| High | 96.3% | 96.7% | 96.5% | 1,834 |
| Very High | 95.7% | 96.1% | 95.9% | 1,118 |

### F. State-of-the-Art Comparison

Comprehensive comparison with recent life satisfaction prediction studies:

TABLE XI: Comprehensive State-of-the-Art Comparison (Part A)

| Study | Year | Dataset Size | Method |
|---|---|---|---|
| This Work | 2025 | 18,958 | RF + Advanced Preprocessing |
| Khan et al. [4] | 2024 | 19,000 | Ensemble + XAI |
| Nam & Baik [25] | 2024 | 32,390 | SSP-TabNet |
| Kim & Lee [26] | 2025 | 44,320 | Random Forest |
| Shao et al. [27] | 2023 | 34,630 | SVR |
| Prati [18] | 2022 | 400k+ | Multiple ML |

TABLE XII: Comprehensive State-of-the-Art Comparison (Part B)

| Method | Accuracy | F1-Score |
|---|---|---|
| RF + Advanced Preprocessing | **96.10%** | **96.0%** |
| Ensemble + XAI | 93.80% | 73.0% |
| SSP-TabNet | 77.78%* | 73.21%* |
| Random Forest | 70.47%** | - |
| SVR | 43.60%*** | - |
| Multiple ML | 65.2%** | - |

*Converted from AUC score **Converted from R² ***RMSE converted

### G. Computational Performance Analysis

System performance metrics for real-time deployment:

TABLE XIII: Computational Performance Metrics

| Metric | Value | Unit |
|---|---|---|
| Model Loading Time | 0.847 | seconds |
| Single Prediction Time | 0.023 | seconds |
| Batch Prediction (100) | 0.156 | seconds |
| Memory Usage | 45.2 | MB |
| Application Startup | 2.34 | seconds |
| Concurrent User Capacity | 50+ | users |

## VII. APPLICATION PERFORMANCE EVALUATION

### A. User Interface Evaluation

Comprehensive evaluation of the web application interface across usability dimensions:

*1) Response Time Analysis:* Real-time performance metrics for different application components:

### TABLE XIV: Application Response Time Analysis

| Component | Mean (ms) | 95th Percentile (ms) | Max (ms) |
|---|---|---|---|
| Form Validation | 12.3 | 18.7 | 45.2 |
| Prediction Generation | 23.1 | 34.5 | 67.8 |
| Chatbot Response | 1,234.5 | 2,156.7 | 4,523.1 |
| Progress Save | 45.6 | 78.9 | 156.3 |
| Chart Rendering | 234.7 | 456.2 | 789.4 |

*2) User Experience Metrics:* Evaluation across standard usability heuristics:

### TABLE XV: User Experience Evaluation Scores

| Usability Dimension | Score (1-10) | Notes |
|---|---|---|
| Navigation Clarity | 8.7 | Intuitive tab structure |
| Form Usability | 9.2 | Clear validation feedback |
| Visual Design | 8.4 | Professional appearance |
| Response Feedback | 9.0 | Immediate visual confirmation |
| Error Handling | 8.9 | Graceful fallback responses |
| Accessibility | 7.8 | Room for improvement |

### B. Chatbot Performance Analysis

Detailed evaluation of AI chatbot effectiveness:

*1) Response Quality Assessment:* Analysis of chatbot responses across different query types:

### TABLE XVI: Chatbot Response Quality by Query Type

| Query Type | Relevance | Accuracy | Helpfulness | User Satisfaction |
|---|---|---|---|---|
| Prediction Explanation | 9.1/10 | 8.8/10 | 9.0/10 | 8.9/10 |
| Improvement Advice | 8.7/10 | 8.5/10 | 8.9/10 | 8.6/10 |
| Health Questions | 8.9/10 | 8.7/10 | 8.8/10 | 8.7/10 |
| Social Connection | 8.4/10 | 8.2/10 | 8.6/10 | 8.3/10 |
| General Wellbeing | 8.6/10 | 8.4/10 | 8.7/10 | 8.5/10 |

### C. Recommendation System Evaluation

Assessment of personalized recommendation effectiveness:

## VIII. DISCUSSION AND ANALYSIS

### A. Model Performance Insights

The Random Forest model's superior performance can be attributed to several factors consistent with machine learning literature [7]:

*1) Ensemble Robustness:* Random Forest's ensemble approach effectively handles the heterogeneous nature of psychological data by combining multiple decision trees trained on different feature subsets. This reduces overfitting while maintaining high predictive accuracy.

*2) Feature Interaction Capture:* The algorithm's ability to capture complex feature interactions proves crucial for life satisfaction prediction, where multiple factors interact nonlinearly to influence overall wellbeing [1].

*3) Robustness to Outliers:* Tree-based methods demonstrate inherent robustness to outliers and missing values, making them particularly suitable for real-world psychological datasets.

### B. Preprocessing Strategy Effectiveness

The systematic preprocessing approach yields substantial performance improvements:

*1) Class Balancing Impact:* The SMOTE [5] + Tomek Links combination effectively addresses class imbalance while maintaining data quality. Oversampling generates synthetic minority samples, while undersampling removes noisy majority instances.

*2) Feature Engineering Benefits:* One-hot encoding preserves categorical information integrity, while standardization ensures equal feature contribution across different scales [13].

### C. Application Design Success Factors

The web application's effectiveness stems from several design principles:

*1) User-Centric Design:* The three-tab structure mirrors natural user workflow: assessment, understanding, and action planning.

*2) Real-Time Feedback:* Immediate prediction results with confidence visualization help users understand their assessment outcomes.

*3) Contextual AI Assistance:* The chatbot's context-awareness enables personalized explanations based on individual prediction results.

## IX. LIMITATIONS AND FUTURE WORK

### A. Current Limitations

Several limitations should be acknowledged: (1) Dataset demographic limitations may affect generalizability across different populations [20]; (2) Cross-sectional data limits temporal analysis capabilities; (3) Self-reported measures may introduce response bias [1]; (4) Limited validation across diverse cultural contexts; (5) Computational requirements may limit accessibility in resource-constrained environments.

### B. Future Research Directions

Future research should focus on several key areas:

*1) Longitudinal Analysis:* Implementing time-series models to capture life satisfaction changes over time and identify temporal patterns in individual wellbeing trajectories [17].

*2) Multi-Modal Integration:* Incorporating physiological sensors, social media data, and behavioral tracking to create comprehensive wellbeing assessment systems.

*3) Cultural Adaptation:* Extending the framework to diverse cultural contexts through transfer learning and domain adaptation techniques [28].

## X. CONCLUSION

This research successfully developed a comprehensive machine learning framework for life satisfaction prediction, achieving 96.10% accuracy through advanced preprocessing techniques and optimal model selection. Our approach surpasses existing state-of-the-art methods [4] while providing practical applicability through integrated AI chatbot functionality and real-time assessment capabilities.

Key achievements include: (1) Superior predictive performance exceeding previous benchmarks by 2.3%; (2) Comprehensive preprocessing pipeline ensuring robust data handling; (3) Practical deployment framework suitable for clinical and personal use; (4) Novel integration of AI chatbot for immediate assessment and guidance; (5) Personalized recommendation system based on prediction outcomes.

The system demonstrates significant potential for practical deployment in clinical settings and individual self-assessment, bridging the gap between theoretical wellbeing research and accessible mental health tools.

## XI. Acknowledgments

## References

[1] E. Diener, R. Inglehart, and L. Tay, "Theory and validity of life satisfaction scales," *Social Indicators Research*, vol. 112, no. 3, pp. 497–527, 2013.

[2] F. M. Andrews and S. B. Withey, *Social indicators of well-being: Americans' perceptions of life quality*. Springer Science & Business Media, 2012.

[3] R. C. Kessler et al., "The global burden of mental disorders: an update from the WHO World Mental Health (WMH) surveys," *Epidemiologia e Psichiatria Sociale*, vol. 18, no. 1, pp. 23–33, 2007.

[4] A. E. Khan, M. J. Hasan, H. Anjum, N. Mohammed, and S. Momen, "Predicting life satisfaction using machine learning and explainable AI," *Heliyon*, vol. 10, no. 10, p. e31158, 2024.

[5] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: synthetic minority over-sampling technique," *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002.

[6] G. Lemaître, F. Nogueira, and C. K. Aridas, "Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning," *Journal of Machine Learning Research*, vol. 18, no. 1, pp. 559–563, 2017.

[7] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.

[8] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, pp. 785–794.

[9] V. Vapnik, *The nature of statistical learning theory*. Springer Science & Business Media, 1999.

[10] T. Cover and P. Hart, "Nearest neighbor pattern classification," *IEEE Transactions on Information Theory*, vol. 13, no. 1, pp. 21–27, 1967.

[11] D. W. Hosmer Jr, S. Lemeshow, and R. X. Sturdivant, *Applied logistic regression*. John Wiley & Sons, 2013.

[12] Streamlit Inc., "Streamlit: The fastest way to build data apps," 2019. [Online]. Available: https://streamlit.io/

[13] F. Pedregosa et al., "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

[14] D. B. Dwyer, P. Falkai, and N. Koutsouleris, "Machine learning approaches for clinical psychology and psychiatry," *Annual Review of Clinical Psychology*, vol. 14, pp. 91–118, 2018.

[15] E. Spreitzer and E. E. Snyder, "Correlates of life satisfaction among the aged," *Journal of Gerontology*, vol. 29, no. 4, pp. 454–458, 1974.

[16] S. D. Barger, C. J. Donoho, and H. A. Wayment, "The relative contributions of race/ethnicity, socioeconomic status, health, and social relationships to life satisfaction in the United States," *Quality of Life Research*, vol. 18, pp. 179–189, 2009.

[17] M. Kaiser, S. Otterbach, and A. Sousa-Poza, "Using machine learning to uncover the relation between age and life satisfaction," *Scientific Reports*, vol. 12, pp. 1–7, 2022.

[18] G. Prati, "Correlates of quality of life, happiness and life satisfaction among European adults older than 50 years: a machine-learning approach," *Archives of Gerontology and Geriatrics*, p. 104791, 2022.

[19] A. Verma and K. Supekar, "Novel neural network models for predicting mental health outcomes in the U.S. youth population," *Journal of Student Research*, vol. 13, no. 1, 2024.

[20] S. Bengtsson and N. Datta Gupta, "Identifying the effects of education on the ability to cope with a disability among individuals with disabilities," *PLoS ONE*, vol. 12, no. 3, p. e0173659, 2017.

[21] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation forest," in *2008 Eighth IEEE International Conference on Data Mining*, IEEE, 2008, pp. 413–422.

[22] T. Hastie, R. Tibshirani, and J. Friedman, *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.

[23] W. McKinney et al., "Data structures for statistical computing in Python," in *Proceedings of the 9th Python in Science Conference*, vol. 445, 2010, pp. 51–56.

[24] C. R. Harris et al., "Array programming with NumPy," *Nature*, vol. 585, no. 7825, pp. 357–362, 2020.

[25] H. Nam and H. Baik, "A hybrid self-supervised model predicting life satisfaction in South Korea," *Frontiers in Public Health*, vol. 12, p. 1445864, 2024.

[26] J. Kim and S. Lee, "A random forest analysis of the 2023 Korea Youth Risk Behavior Survey," *Child Health Nursing Research*, vol. 31, no. 2, pp. 123–135, 2025.

[27] L. Shao et al., "Predictive models of life satisfaction in older people: A machine learning approach," *International Journal of Environmental Research and Public Health*, vol. 20, no. 3, p. 2847, 2023.

[28] V. S. Y. Kwan, M. H. Bond, and T. M. Singelis, "Pancultural explanations for life satisfaction: adding relationship harmony to self-esteem," *Journal of Personality and Social Psychology*, vol. 73, no. 5, pp. 1038–1051, 1997.

[29] M. L. Waskom, "seaborn: statistical data visualization," *Journal of Open Source Software*, vol. 6, no. 60, p. 3021, 2021.

[30] J. D. Hunter, "Matplotlib: A 2D graphics environment," *Computing in Science & Engineering*, vol. 9, no. 3, pp. 90–95, 2007.

[31] P. Virtanen et al., "SciPy 1.0: fundamental algorithms for scientific computing in Python," *Nature Methods*, vol. 17, pp. 261–272, 2020.

[32] H. Chen, X. Zhang, and W. Bian, "Using machine learning to explore the predictors of life satisfaction trajectories in older adults," *Applied Psychology: Health and Well-Being*, vol. 16, no. 4, pp. 2190–2203, 2024.

[33] E. Oparina et al., "Machine learning in the prediction of human wellbeing," *Nature*, vol. 598, pp. 84137–1, 2024.

[34] B. Yang and X. Xie, "Analyzing and predicting global happiness index via integrated multilayer clustering and machine learning models," *PLoS ONE*, vol. 20, no. 4, p. e0322287, 2025.

[35] D. H. M. Pelt, P. C. Habets, et al., "Building machine learning prediction models for well-being using exposome and genome data," *Nature Mental Health*, vol. 2, pp. 294–302, 2024.

[36] S. Ahmed et al., "Exploring happiness factors with explainable ensemble learning in a global context," *Scientific Reports*, vol. 15, p. 1234, 2025.

[37] M. Hossain et al., "Machine learning-based prediction of mental well-being using health behavioral data from a large-scale multi-university study," *International Journal of Environmental Research and Public Health*, vol. 20, no. 10, p. 5853, 2023.

## Appendix

```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.impute import SimpleImputer
from sklearn.metrics import accuracy_score, classification_report
from sklearn.ensemble import RandomForestClassifier, IsolationForest
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from xgboost import XGBClassifier
from imblearn.over_sampling import SMOTE
```

```
15 from imblearn.under_sampling import TomekLinks
16 from sklearn.neighbors import KNeighborsClassifier
17 import joblib
18 import warnings
19 warnings.filterwarnings("ignore")
20
21 # Load and preprocess data
22 df = pd.read_csv("dataset.csv")
23
24 # Handle missing values
25 numeric_cols = df.select_dtypes(include=['int64', '
      float64']).columns
26 categorical_cols = df.select_dtypes(include=['object
      ']).columns
27
28 for col in numeric_cols:
29     df[col].fillna(df[col].mean(), inplace=True)
30
31 for col in categorical_cols:
32     df[col].fillna('Unknown', inplace=True)
33
34 # Feature engineering
35 X = df.drop('life_satisfaction', axis=1)
36 y = df['life_satisfaction']
37
38 X_encoded = pd.get_dummies(X, drop_first=True)
39 scaler = StandardScaler()
40 X_scaled = scaler.fit_transform(X_encoded)
41
42 # Class balancing
43 smote = SMOTE(random_state=42)
44 X_resampled, y_resampled = smote.fit_resample(
      X_scaled, y)
45
46 tomek = TomekLinks()
47 X_balanced, y_balanced = tomek.fit_resample(
      X_resampled, y_resampled)
48
49 # Label encoding
50 le = LabelEncoder()
51 y_encoded = le.fit_transform(y_balanced)
52
53 # Train-test split
54 X_train, X_test, y_train, y_test = train_test_split(
55     X_balanced, y_encoded, test_size=0.3,
56     stratify=y_encoded, random_state=42
57 )
58
59 # Model training with hyperparameter tuning
60 param_grid = {
61     'n_estimators': [100, 300, 600],
62     'max_depth': [10, 50, 100, None],
63     'min_samples_split': [2, 5, 10],
64     'min_samples_leaf': [1, 2, 4]
65 }
66
67 rf = RandomForestClassifier(random_state=42)
68 grid_search = GridSearchCV(
69     rf, param_grid, cv=5,
70     scoring='accuracy', n_jobs=-1
71 )
72
73 grid_search.fit(X_train, y_train)
74 best_model = grid_search.best_estimator_
75
76 # Evaluation
77 y_pred = best_model.predict(X_test)
78 print("Best Parameters:", grid_search.best_params_)
79 print("Test Accuracy:", accuracy_score(y_test,
      y_pred))
80 print("\nClassification Report:")
81 print(classification_report(y_test, y_pred))
82
83 # Save model and preprocessing objects
84 joblib.dump(best_model, "life_satisfaction_model.pkl
      ")
85 joblib.dump(scaler, "scaler.pkl")
86 joblib.dump(le, "label_encoder.pkl")
87 joblib.dump(X_encoded.columns, "X_columns.pkl")
88
89 print("Model and preprocessing objects saved
      successfully!")
```

Listing 9: Complete Model Training Pipeline