# Explainable Artificial Intelligence for Identifying Psychological Risk Profiles of Youth Suicidal Ideation: A SHAP-Based Machine Learning Analysis

Mariana. Coutinho[1]*, Karim. Fahmy[2]

[1] Department of Clinical Psychology, Federal University of Rio de Janeiro, Rio de Janeiro, Brazil
[2] Department of Clinical Psychology, Cairo University, Giza, Egypt

* Corresponding author email address: mariana.coutinho@ufrj.br

| Editor | Reviewers |
|---|---|
| Trevor Archer | **Reviewer 1:** Zahra Yousefi |
| Professor Department of Psychology University of Gothenburg, Sweden trevorcsarcher49@gmail.com | Assistant Professor, Department of Psychology, Khorasgan Branch, Islamic Azad University, Isfahan, Iran. Email: yousefi1393@khuisf.ac.ir |
|  | **Reviewer 2:** Mehdi Rostami |
|  | Department of Psychology and Counseling, KMAN Research Institute, Richmond Hill, Ontario, Canada. Email: dr.mrostami@kmanresce.ca |

## 1. Round 1

### 1.1. Reviewer 1

Reviewer:

In the statement "Youth suicidal ideation is not a unitary phenomenon but rather the result of complex interactions…", the manuscript introduces a key conceptual premise. However, the paragraph would benefit from an explicit articulation of how this complexity motivates the use of machine learning rather than traditional multivariate models. Consider adding one bridging sentence clarifying this methodological necessity.

The discussion of black-box concerns is conceptually strong; however, the manuscript would benefit from a clearer ethical framing. Consider explicitly stating why lack of interpretability is especially problematic in suicide risk contexts compared to other prediction domains (e.g., false positives, clinical accountability).

The final sentence outlining the study aim is clear, but it would be strengthened by explicitly stating the primary outcome type (binary vs. continuous suicidal ideation) and the intended level of application (research, clinical decision support, school-based screening).

The phrase "a sufficiently large cohort to support machine learning model training" is vague. Please report an explicit rationale for sample adequacy (e.g., events-per-variable logic, class balance justification, or reference to ML sample size heuristics).

You note recruitment from "urban and semi-urban regions of central and southern Chile." Please clarify whether regional clustering effects were examined or adjusted for, as this may introduce contextual dependence in both predictors and outcomes.

The manuscript states that multiple models were trained but does not clarify whether feature selection was embedded or external. Please specify whether all predictors entered each model simultaneously or whether dimensionality reduction or regularization was applied.

Authors uploaded the revised manuscript.

*1.2.    Reviewer 2*

Reviewer:

The paragraph beginning with "Certain youth populations experience disproportionate risk…" lists several high-risk groups. Please clarify whether these groups were explicitly represented, measured, or controlled for in your sample. Otherwise, the paragraph risks implying empirical coverage that the study may not substantively provide.

The claim "machine learning methods have demonstrated superior predictive performance" is accurate but underspecified. Please briefly indicate why ensemble methods are particularly suitable for suicide risk modeling (e.g., non-linearity, interaction effects, collinearity tolerance), ideally linking this to the variables included in your dataset.

The manuscript states that suicidal ideation was assessed using "a validated measure" but does not name the instrument. For transparency and reproducibility, the exact scale name, number of items, scoring method, and clinical cutoff (if applicable) should be specified.

Several constructs are listed (e.g., emotional dysregulation, impulsivity, problematic digital use), yet internal consistency indices (e.g., Cronbach's alpha or omega) are not reported. Please include reliability estimates for the current sample, not only prior validations.

The handling of missing data is described as "appropriate imputation techniques" without specification. Given the sensitivity of suicide research, please explicitly report the imputation method (e.g., median imputation, MICE) and justify its suitability for the data structure.

Authors uploaded the revised manuscript.

## 2.    Revised

Editor's decision after revisions: Accepted.
Editor in Chief's decision: Accepted.