

ALL MODELS ARE BIASED, SOME ARE MORE TRANSPARENT ABOUT IT: FULLY INTERPRETABLE AND ADJUSTABLE MODEL FOR MENTAL DISORDER DIAGNOSIS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent advances in machine learning have enabled AI applications in mental disorder diagnosis, but many methods remain black-box or rely on post-hoc explanations which are not straightforward or actionable for mental health practitioners. Meanwhile, interpretable methods, such as k-nearest neighbors (k-NN) classification, struggle with complex or high-dimensional data. A network-based k-NN model (NN-kNN) combines the interpretability with the predictive power of neural networks. The model prediction can be fully explained in terms of activated features and neighboring cases. We experimented with the model to predict the risks of depression and interviewed practitioners. The feedback of the practitioners emphasized the model’s adaptability, integration of clinical expertise, and transparency in the diagnostic process, highlighting its potential to ethically improve the diagnostic precision and confidence of the practitioner.

1 INTRODUCTION

The booming era of Artificial Intelligence sees a rise of its application to mental health-related issues (Graham et al., 2019) and specifically to mental disorder diagnosis (“diagnosis” for short) (Bzdok & Meyer-Lindenberg, 2018; Chattopadhyay, 2017; Graham et al., 2019; Iyortsuun et al., 2023). Although various AI technologies have achieved high accuracies in diagnosis, they generally lack explanation for their decision making, and the process remains a black box for both mental health practitioners (“practitioners” for short) and clients seeking counseling services (“clients” for short) (Jarvie & Lindén, 2024; Lau et al., 2023). Therefore, there is still much caution and concern about the use of AI for diagnosis (Kerz et al., 2023; Li et al., 2024).

This gives rise to recent trends in using Explainable AI (XAI) technologies for mental health: Practitioners and clients may rely on AI tools “to the extent they can economise on human oversight, monitoring and verification of the system’s outputs” (Joyce et al., 2023). Most XAI methods are based on post-hoc explanation methods such as SHAP (Shapley, 1953) and LIME (Ribeiro et al., 2016). However, Rudin (2019) argues that post-hoc explanations are often inadequate for black-box models. The explanations may justify the model’s decision after the decision is made, but not how the model reached the decision internally. In the worst cases, post-hoc explanations are excuses for a model’s mistakes and offer no opportunity to debug or fine-tune a trained black-box model.

To address the interpretability and adjustability of AI models that facilitate practitioners in diagnosis, we propose to use a recently invented model, a neural network based k-nearest-neighbor algorithm (NN-kNN). As a k-nearest neighbor algorithm, NN-kNN can explain each model decision with activated cases, and each activated case can be attributed to its feature distances with the query. As a neural network, NN-kNN allows end-to-end training for both feature weights and case weights, and is compatible with other neural network methods (Ye et al., 2024).

Our study introduces a novel approach to human-machine interaction drawing on insights and methodology from both AI and psychology. Specifically, we study how NN-kNN facilitates practitioners in aspects beyond traditional diagnosis, including feature exploration, explanation of decisions, past case retrieval, and manual weight adjustment. We conducted a qualitative study, using

interpretative phenomenological methods to capture the insights of the practitioner’s experience with the interpretable and adjustable model. Although qualitative research is common in the field of psychology, incorporating this method in AI research offers a novel perspective. This approach enriches existing XAI research and expands the understanding of human interaction with machine learning models, as discussed in Section 4.3.

This article is organized as follows. Section 2 describes the interpretable model we use. Section 3 discusses related work in both the fields of explainable AI and mental health diagnosis. Sections 4 and 5 describe the qualitative study method and findings. Lastly, the article concludes with discussions and future directions.

2 NEURAL NETWORK BASED K-NEAREST NEIGHBORS

2.1 K-NEAREST NEIGHBORS CLASSIFIERS

K-nearest neighbors classifiers (k-NN) examine the task domains involving cases C in the form of (x, L_x) . For each case $x \in C$, $f(x) = \langle x_1, x_2, \dots, x_m \rangle$ is a feature vector describing a client’s information (survey answers, medical record information, etc.) and L_x is the diagnosis label associated with the client (risk of depression). The function f might describe a feature extraction method or simply use the surface features of x . A naive k-NN calculates the distance between two cases as the sum of distances between corresponding features, using a simple Minkowski distance measure such as the Euclidean distance (Dasarathy, 1991). Better k-NN methods use feature weights in the calculation of distance. Traditional methods use a global feature weighting, while others allow certain sets of cases (or even each individual case) to have their own feature weighting (Manzali et al., 2024; Aha & Goldstone, 1992; Friedman, 1994; Ricci & Avesani, 1995; Marchiori, 2013; Bonzano et al., 1997). Some methods assign case weights, so that certain cases contribute more in the voting of the final prediction (Bicego & Loog, 2016; Aguilera et al., 2019). The case weights may be based on the distance between the case and the query, or it may be a learned parameter of the case. Other methods such as neighborhood components analysis (NCA) and large margin nearest neighbor (LMNN) transform the feature space to extract high-level features before distance calculation (Goldberger et al., 2004; Weinberger & Saul, 2009).

2.2 NEURAL NETWORK BASED K-NEAREST NEIGHBORS ALGORITHM

The neural network based k-nearest neighbors algorithm (NN-kNN) implements both feature weights and case weights by having network layers that simulate the behavior of a k-NN:

1. The case layer stores all cases (the training set).
2. The feature extraction layer extracts the features of the query q and each case x .

$$f(q) = \langle q_1, q_2, \dots, q_m \rangle, f(x) = \langle x_1, x_2, \dots, x_m \rangle \quad (1)$$

3. The feature distance layer calculates the distance between the corresponding features as

$$\delta_i = \delta_i(q_i, x_i) \geq 0 \quad (2)$$

We choose the squared distance $\delta_i(q_i, x_i) = (q_i - x_i)^2$

4. The case activation layer activates a case x given the query q by

$$case_activation(x|q) = \sigma(w_{x\delta_1} * \delta_1 + w_{x\delta_2} * \delta_2 + \dots + b_x) \quad (3)$$

where $w_{x\delta_i}$ is the weighting of feature i for case x and $b_x \geq 0$ is the default activation for case x . $w_{x\delta_i}$ can only be negative because the more different the case from the query, the less activated the case. We choose the sigmoid function $\sigma()$ as the activation function to limit *case_activation* within $[0, 1]$.

5. The top k case selection layer is optional. When enabled, it will keep the top k case activations and resets other case activations to 0.
6. The target activation layer takes each case’s activation to activate the corresponding class.

$$class_activation(L|q) = \sum_x (w_{(x,L)} * case_activation(x)) + b_L \quad (4)$$

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

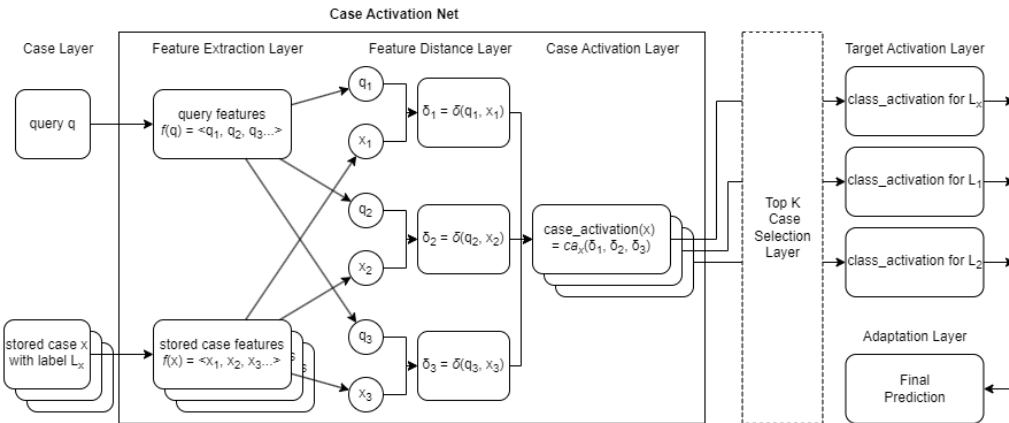


Figure 1: The Model of NN-kNN (Ye et al., 2024)

where $w_{(x,L)}$ is the weight of the case x for the class L and b_L is the bias of the class L . $w_{(x,L)}$ is forced to be positive (by using a ReLU) if the case x is of class L and $w_{(x,L)} = 0$ otherwise, because a case should only activate its corresponding class but not other classes.

7. The adaptation layer chooses the label L with the maximum *class_activation*.

The layers are depicted in Figure 1. In previous experiments comparing NN-kNN with neural networks and state-of-the-art k-NN methods (such as LMNN), NN-kNN achieves equal or less prediction error in classification and regression on multiple datasets (Ye et al., 2024).

NN-kNN has the potential to incorporate more advanced neural architectures. For instance, additional layers can be added in the feature extraction layer to transform the feature space like NCA or LMNN does. Although adding deeper layers to NN-kNN can enhance its analytical capabilities, it comes at the cost of reduced interpretability. Greater analytical power is not always essential, and interpretable models often achieve good enough performance, sometimes even surpassing non-interpretable models (Rudin, 2019; Ye et al., 2024).

3 RELATED WORK

3.1 EXPLAINABLE AND INTERPRETABLE AI

According to the survey by Joyce et al. (2023), the majority of XAI methods in mental health are based on feature importance methods such as SHAP (Shapley, 1953). Feature importance methods estimate the weights of features by measuring the influence on the model’s output after perturbing a feature. Only two methods in the survey are regression-based and interpretable by design. Similarly, most studies in the wider field of XAI are post-hoc methods because they are flexible and easily applicable to different models (Saleem et al., 2022). Post-hoc methods build an interpretable model that mimics the behavior of a black-box model and use this new model as an explanation. However, post-hoc explanations are problematic for high-stakes decision making in mental health applications: (1) post-hoc explanations may not faithfully represent the original model’s computations, (2) they may lack the detail necessary to fully understand the black-box model, and (3) they do not permit manual calibration by domain experts (Rudin, 2019).

Our XAI method is interpretable by design and explains decisions using features and cases. This is similar to the explanation done in case-based reasoning (CBR) (Schoenborn et al., 2021; Gates & Leake, 2021), where decisions on queries are explained by past cases similar to queries. Most CBR systems weight features (Wettschereck et al., 1997), while fewer weight cases (Bicego & Loog, 2016). NN-kNN takes the extra step to train both feature weights and case weights at the same time in an end-to-end manner.

While most modern machine learning methods learn parameters through training, some interpretable AI systems allow parameters to be set using expert knowledge or external methods. For example, in

k-NN, feature or case weights can be determined using mutual information (García-Laencina et al., 2009), the analytic hierarchy process (Bhattacharya et al., 2017), or fuzzy membership functions (Biswas et al., 2018). NN-kNN offers a hybrid approach: feature and case weights are initially learned by the network, then experts can fine-tune them for retraining.

3.2 MEMORY AUGMENTED NETWORKS

Weston et al. (2014) proposed the class of memory networks. A memory network consists of a memory m and four components: an input feature map I that extracts features of the query, a generalization process G that updates memory given the query, an output feature map O that produces the output using the query and the memory, and lastly a response R converts output to desired format. They implemented a memory network for question answering purpose. In their example, the query is a textual question, m stores texts, O finds k supporting memories and produces an output feature, and R produces a textual response using RNN on the output feature.

Matching networks (Vinyals et al., 2016) extend the memory networks for one-shot learning. The authors do so by incorporating characteristics from non-parametric models, allowing a trained network to be directly used on a new support set. Similarly, prototypical networks (Snell et al., 2017) learn a metric space and perform classification by computing the distances between the query and prototypes of each class. Li et al. (2018) propose a neural network model that stores auto-encoded embeddings of learned prototypes and makes prediction by comparing the query embedding with the prototype embeddings.

NN-kNN fits the class of memory networks: m is the case layer, I is the feature extraction layer, G is the training of network parameters, and O is the prediction layer based on activated cases. NN-kNN is similar to the matching networks and prototypical networks because the layers until the target activation layer serve as a similarity metric for the cases. In fact, NN-kNN can be considered as a generalization of the previous methods. It works with any case (not just prototypes) or any feature (not just embeddings extracted by an autoencoder). Its k-NN nature allows easy insertion and deletion of cases and easy explanation through activated cases and features.

3.3 MENTAL HEALTH DIAGNOSIS BY PRACTITIONERS

The development of efficient diagnosis has stagnated for decades, facing several challenges:

1. **Time-Consuming Diagnostic Processes:** Clinicians continue to rely on diagnostic manuals such as the DSM-V and ICD-10. The time-consuming diagnostic process takes away valuable time from direct therapeutic interventions (Perkins et al., 2018)
2. **Insufficient Training Opportunities:** Only 23% of American Psychological Association (APA)-accredited doctoral programs provide trainee clinicians with the training sites necessary for structured diagnostic interviews (Mihura et al., 2017).
3. **Inadequacy of Diagnostic Manuals:** A systematic review involving 2,228 participants found that clinicians often view diagnostic manuals as unhelpful due to incomplete or inaccurate symptom descriptors (Perkins et al., 2018). They tend to focus on categorizing symptoms rather than identifying the underlying formations of psychological symptoms.
4. **High Rates of Misdiagnosis:** In a study of 309 psychiatric patients, 39.16% patients with severe psychiatric disorders were misdiagnosed, with schizoaffective disorder having the highest misdiagnosis rate (75%), followed by major depressive disorder (54.72%) (Ayano et al., 2021).
5. **Comorbidity and Diagnostic Complexity:** Mental health diagnoses are further complicated by comorbid conditions. For instance, individuals with autism spectrum conditions are more likely to be diagnosed with mental health disorders than non-autistic individuals, increasing the potential for misdiagnosis (Au-Yeung et al., 2019).

Given these ongoing challenges, the field of health service psychology urgently requires a more effective and accurate diagnostic system to support clinicians in making better informed decisions.

3.4 MENTAL HEALTH DIAGNOSIS BY AI

Many studies have demonstrated the efficacy of AI-enhanced tools in supporting clinical diagnosis Graham et al. (2019). For example, Zhang et al. (2022) conducted a comprehensive narrative review of 399 studies on NLP applications in mental illness detection over the past decades. Their finding suggested a significant upward trend in research focused on NLP for mental illness detection, with deep learning approaches showing better performance than traditional machine learning methods. In a recent study, Lau et al. (2023) examined deep learning models to automate depression severity assessment by using parameter efficient tuning techniques and pre-trained large language models. Their results suggested that prefix tuning allowed for more efficient model training to reduce overfitting and offered a scalable solution of automatic depression assessment.

Despite the great promise of AI-assisted clinical diagnosis, both patients and clinicians remain hesitant to accept its application. In one study, patients have expressed a preference for AI tools that are tailored, person-centered, and adaptable to their individual treatment plans and expressed concern that they must adapt to technology (Li et al., 2024). In addition, clinicians often hesitate to fully integrate AI technology into their practice due to concern about their understanding of how AI systems work (Kerz et al., 2023). As such, there is a need for transparency and explainability in AI systems to foster trust with mental health practitioners. However, among the mental health XAI projects surveyed by Joyce et al. (2023), none interviewed specialists about their experience with the XAI models. Surveys in the broader field of XAI (Saeed & Omlin, 2023; Das & Rad, 2020; Molnar et al., 2020) also identified evaluating the explainability of XAI models as a major challenge. In general, there is no agreed upon measure of explainability/interpretability.

In this study, we embrace the idea that explainability is not a rigorous formal concept and that an explanation is only as good as the intended audience perceives it. We take a human-centered approach to evaluate explainability in a qualitative approach that is more common in social science.

4 QUALITATIVE INTERVIEW AND EVALUATION

Previous work has already examined the prediction performance of NN-kNN (Ye et al., 2024). This study focuses on the interpretability and adaptability of NN-kNN to aid practitioners in mental disorder diagnosis. We trained NN-kNN on a dataset on the prediction of depression risk and conducted a qualitative study interviewing 10 licensed practitioners about their experience with the model. Meta-parameters for the computational experiment are described in Appendix B.

4.1 THE DATASET

The model is trained on the dataset from Orozco-del Castillo et al. (2021). The original dataset is answers to a survey of 117 items of true/false answers (see Appendix C). The survey was designed for depression screening according to the Diagnostic and Statistical Manual of Mental Disorders (American Psychiatric Association, 2013) and answered by 157 undergraduate students. After data preprocessing (see Appendix D), our final dataset contains 117 cases; Each case has 102 features of binary values and a class label of 0, 1 or 2, representing low, medium and high risk.

In these experiments, NN-kNN achieves an average accuracy of 0.646, outperforming both standard k-NN (0.417) and LMNN (0.492). We urge caution due to the small dataset size and the unstable results. We did not test extensively on multiple datasets against various models, as previous research (Ye et al., 2024) has already conducted such comparisons. The primary focus of this study is on the interpretability and adjustability of the model for practical use by clinicians.

4.2 THE INTERPRETABILITY AND ADJUSTABILITY FOR PRACTITIONERS

NN-kNN is both interpretable and transparent according to the Transparency and Interpretability for Understandability framework (Joyce et al., 2023). Specifically, when predicting depression risk, the model’s design emphasizes the following key aspects:

- All cases share the same feature weights ($w_{x\delta_i} = w_{y\delta_i} = w_{\delta_i}$ for any two cases x and y), reducing the overall number of parameters. This global feature weighting approach assigns a single weight to each feature across all cases.

- Initial settings for each case x include feature weights set to $w_{x\delta_i} = 1$, case bias at $b_x = 50$, case weight at $w_{(x,L)} = 1$, and class biases at $b_L = 1$. In this configuration, all features and cases start with equal importance. Through training, the model adjusts these weights to improve prediction accuracy. Practitioners can then examine the resulting weights to determine which features or cases the model has prioritized.
- The calculations involved in feature distances, case activations, and class activations rely on basic operations like subtractions and summations. This simplicity ensures that practitioners can easily understand the model’s operations. In addition, they can choose to focus on the most or least weighted features/cases for a quick overview of the model’s results.
- Each parameter in the model has an interpretable role. The feature weight $w_{x\delta_i}$ and the case weight $w_{(x,L)}$ respectively reflect the relevance of a feature/case to depression. The case bias b_x and the class bias b_L respectively indicate the inherent importance of a case and a class. Practitioners can easily identify outliers or particularly influential features/cases.
- Because each parameter has a specific semantic meaning and a preset initial value, practitioners can manually adjust the feature or case weights based on their expert knowledge. After making adjustments, the model can retrain to further refine the predictions. Practitioners can then assess whether the updated model aligns better with their clinical expertise. We demonstrated this functionality to practitioners and collected their feedback on the adjustable model.

4.3 QUALITATIVE INTERVIEW DESIGN

The effectiveness of an explanation ultimately depends on how well it is understood and perceived by its intended audience. To better understand the experiences and perceptions of mental health practitioners, we conducted interviews with 10 licensed clinicians from the US using the interpretative phenomenological analysis (IPA) approach, a widely recognized qualitative research method (Eatough & Smith, 2017). While quantitative methods are prevalent in XAI research—often involving comparisons of models across datasets and using metrics such as prediction accuracy—we chose a qualitative approach for several key reasons:

- **Focusing on a specialized population:** Our target audience consists of licensed clinicians, a small and specialized group. A qualitative approach allows us to zoom in on this specific population, gaining in-depth insights into their experiences with the model, which could be missed in a broader, quantitative study.
- **Tailored model demonstrations:** Demonstrating our model to potential users individually allows us to observe first-hand how each clinician interacts with the system. This personalized approach helps capture the intricacies of their responses, including any challenges they face or specific features they find most valuable.
- **Understanding nuanced reactions:** AI models can evoke a range of responses, particularly when introduced in sensitive fields such as mental health. A qualitative approach is well suited to capturing these nuanced reactions, such as concerns about integrating AI into clinical practice, the potential for AI to aid in diagnostic processes, or hesitations about interpretability. This method allows for a deeper understanding of users’ trust and confidence in the model (Maxwell, 2021).
- **Building trust and bridging theory to practice:** By engaging directly with clinicians through interviews, we can address their concerns, answer questions, and foster trust in the model’s practical use. This is a crucial step in translating theoretical AI advancements into real-world applications that mental health practitioners are willing to adopt.

Four licensed psychologists at the doctorate level and five licensed clinicians at the master level participated in the interviews, among which four identified as men and six as women. Individually, we demonstrated our model in a Jupyter notebook through Zoom, then invited practitioners to adjust the parameters based on their clinician judgements for depression, followed by 30 minutes devoted for them to answer our qualitative questions. Sample interview questions include: “did you feel like the model has become more useful clinically after you tuned the feature weight?” and “since our model can detect bias, if the model’s explanations differ from your clinical judgements, what would you do?” (See Appendix E for a full list of questions).

4.4 DATA ANALYSIS THROUGH INTERPRETATIVE PHENOMENOLOGICAL METHOD

Following the four-step analytic process of Interpretative Phenomenological Analysis (IPA), a team of three members (a licensed counseling psychologist, a doctoral candidate in counseling psychology, and an undergraduate psychology student) began by checking each other's biases related to AI and clinical diagnosis. Each member then independently reviewed the interview transcripts, annotating their initial reactions. These annotations were translated into experiential statements, summarizing key aspects of each participant's experience. Next, we compared our individual statements to resolve any discrepancies and clustered them to generate overarching themes. Finally, we compiled eight themes that broadly reflected most participants' experiences. Before finalizing the results, the third author acted as an auditor, reviewing that the themes were grounded in the original transcripts and participants' experiences.

5 INTERVIEW FINDINGS

According to our data analysis, eight themes were generated, endorsed by most (at least 6) if not all participants. To remain consistent with the typical way of reporting results for IPA studies (Liu et al., 2020), we combined the themes with our interpretations of participants' experiences to demonstrate a contextual exploration of licensed clinicians' perceptions. To ensure confidentiality, all names used in the following sections are pseudonyms.

5.1 THEMES IDENTIFIED AND DISCUSSIONS

Theme I: Building Trust and Precision in Diagnosis Participants initially perceived AI-generated diagnoses as inaccurate and unusable, but changed their view after engaging with the adjustable model. After attending our demonstration and adjusting the models themselves, participants reported an appreciation of this model's ability to tune feature weights and adjust cases, allowing clinicians to focus on the most clinically relevant aspects and make diagnoses more precise (Dr. Nate):

“The tuning process did help somewhat in clarifying the approach to diagnosing by adjusting the feature weights. It allowed for a more targeted focus on clinically relevant aspects.”

Participants also emphasized the importance of transparently presenting the model's diagnostic process, which would enhance clinicians' confidence in considering AI-generated results (Dr. Yun):

“Given that the model is already explainable, and I can see the whole process regarding coming to the diagnosis. It feels more comfortable to understand how it comes to the conclusion. If I want to make some changes to certain features, I can also do that. This gives me more confidence in terms of using it in clinical situations.”

Besides, participants observed the adjustment process increased accuracy of depression diagnosis:

“A general screening can sometimes overlook important clinical features of a client. By being able to adjust the weights on specific features that impact the client's symptoms, I believe the screening becomes more clinically accurate and personalized to the individual, which enhances the overall assessment and treatment planning process.”

Theme II: Potential Risk of Bias Introduced by Model Adjustment Most participants identified the challenge of maintaining a balance between flexibility and accuracy, expressing concern that their adjustments might lead to less accurate diagnoses (Xing):

“I am worried about the weight of some questions as I was tuning, which are not the signature of depression...I didn't not quite understand the decision making process of the AI because I would disagree with its clinical decision.”

Participants also voiced concerns about the potential risk of introducing bias through adjustments, as clinicians might unintentionally influence the model to confirm pre-existing beliefs (Dr. Yong):

“The ability to tune the model increases the risk of introducing bias or misusing the tool, especially if there isn't enough transparency about how tuning decisions are made and under what circumstances. Without clear guidelines and a full understanding of the impact of tuning, the potential for bias could compromise trust rather than enhance it.”

378 **Theme III: Customization for Clinical Expertise and Multicultural Factors** Several participants
379 appreciated the ability to adjust feature weights. They emphasized the model’s flexibility to ac-
380 commodate their clinical experiences and theoretical preferences regarding diagnoses (Dr. Nate and
381 Lily):

382 “The tunable feature allows for more tailored diagnostic criteria, which can be useful in
383 specific cases, settings, and with particular clients. It offers flexibility, helping clinicians
384 adjust for individual circumstances.”

385 “It is customized and flexible to fit the patient population (anxiety in kids with or without
386 chronic illness may look different)”

387
388 Many clinicians also noted the model’s ability to account for multicultural factors, allowing them to
389 adjust for diverse clients and incorporate cultural nuances into the diagnostic process (Dr. Yun):

390 “After I tuned the feature weights, the AI model will run again based on my input and gen-
391 erate new features/cases that meets the standards, which is very intelligent. Additionally,
392 regarding the multicultural consideration, clinicians can also tune the features related to
393 clients from diversity background. This is pretty cool.”

394 Furthermore, Dr. Yalin expressed that the model feels like a reliable assistant or skilled trainee, who
395 undergoes personalized training from the clinician to ultimately save time and double-check clinical
396 decisions:

397 “I felt like I was training a reliable trainee that, once training completed, can be a great time
398 saver and double check my clinical judgement.”

400 **Theme IV: Transparency and Ethical Trustworthiness** Participants shared their perspectives on
401 the transparency and ethical trustworthiness of the model, which were consistently identified as key
402 factors influencing their willingness to use it in clinical practice. They highlighted the model’s en-
403 hanced transparency, noting how it distinguishes itself from other AI technologies (Lina and Yaya):

404 “It can be ethically reliable because it ensures informed consent from clients and provides
405 transparency in how it influences counselors’ decision-making.”

406 “Additionally, transparency in how AI generates recommendations can further build trust.
407 Ultimately, while trust varies among clinicians, the adaptability of tunable AI can signifi-
408 cantly improve its integration into patient care.”

409 However, many clinicians raised ethical concerns about potential misuses, especially if clinicians
410 lack clear clinical or technical guidelines. Similarly, Dr. Yong expressed ongoing concerns about
411 whether the model is truly more ethically trustworthy, despite its technical advancements. She
412 also voiced apprehension about potential confusion regarding the role of AI in replacing human
413 judgment:

414 “um...I would say it’s better than ChatGPT. I am not sure about ‘more ethically trustworthy.’
415 There are still unclear areas seem to be addressed.”

416 “I think it would be helpful to state that the model is not aiming to take away the human
417 factors when introducing the algorithm. There were moments I got confused from last and
418 this time that the algorithm is trying to take over the human factors in the clinical decision
419 making.”

420 **Theme V: Differential Diagnosis and Time Efficiency** Participants shared other unique strengths
421 of our AI model, especially when it comes to aiding certain differential diagnosis as well as saving
422 time clinically. Before trying out our model, participants had concerns around the depths of AI-
423 generated diagnoses, given that their experiences using generative AI, such as ChatGPT has been
424 frustrating and superficial. Thus, our explainable AI model impressed them with its transparency,
425 resulting in deeper reflections around specific ways of using it for clinical diagnosis.

426 For example, Dr. Yun shared that in her opinions, AI could assist in quicker and more accurate dif-
427 ferential diagnoses, such as distinguishing between Major Depressive Disorder (MDD) and Paranoid
428 Personality Disorder (PPD):

429
430 “Another feature could be helpful is the differential diagnosis, such as MDD/ PPD. If AI
431 can show why it’s MDD instead of PDD, that it will certainly help the clinician to spend
less time in terms of making diagnosis.”

432 Dr. Yalin disclosed that the traceability and documentation of our model’s steps make it a valu-
433 able resource for complex decision-making, especially around ethical issues and when dealing with
434 differential diagnoses:

435 “because each step can be traced and documented. In our hospital, we often have ethics
436 board meeting to review complex cases (mostly in the medical side) but I can see how
437 having the algorithm can be a supportive data to aid with decision making.”

438 “Not sure if this is doable - perhaps having a model that can pick up on comorbidity and dif-
439 ferential diagnoses? For example, it is sometimes hard to differentiate between childhood
440 anxiety, ADHD, and kids who have both disorders.”

441 Other participants believed that with refinement, our model has the potential to simplify complex
442 tasks and improve diagnostic efficiency, highlighting the excitement around using AI to streamline
443 clinical processes:

444 “The potential for simplifying complex tasks and improving diagnostic efficiency is excit-
445 ing...I look forward to seeing how it develops further and hope that with some improve-
446 ments, it can greatly assist clinicians.”

448 **Theme VI: Psychotherapy Insights for Future Directions** Along with the excitement of using our
449 AI model for differential diagnosis, some clinicians offered expectations of its future clinical utility.
450 Participants described the ability to track mood symptoms over time offers valuable insights during
451 ongoing therapy, helping clinicians adjust treatment plans based on evolving patient data (Lina):

452 “Also, it is beneficial for tracking mood symptoms over time with more thorough contex-
453 tual insights. I would review the model’s explanations and compare them to my clinical
454 insights. Ongoing psychotherapy will allow counselors to gain a deeper understanding of
455 their clients, enabling them to enhance their clinical judgment and make more informed
456 decisions moving forward.”

457 Other participants thought that AI could enhance diagnostic reasoning by flagging symptoms that
458 don’t fit a particular diagnosis. Dr. Yun expressed that AI can not only show why it supports a
459 diagnosis but also highlight why it rules out others like anxiety or obsessive-compulsive disorder:

460 “I was thinking if the AI can also add some features, such as flagging some symptoms
461 that’s not aligning with diagnosis (i.e. feeling very energetic in everyday life - which is not
462 typical for a depression diagnosis), this could also help the clinician see why the AI make
463 this diagnosis. It’s like not only showing the key symptoms that make AI decide that this
464 case fits diagnosis, but also symptoms that makes AI decide it’s not anxiety/OCD/eating
465 disorder, etc. This could also be helpful for the clinicians to see why AI rule out these other
466 diagnosis.”

467 Last but not least, one participant (Xing) suggest the model to capture key mental health factors, like
468 affect and history, and summarize the client’s mental health background:

469 “I think there are so many important factors that the model is not detecting yet. eg. affects,
470 history of mental health. I would also be curious of a feature that could summarize client’s
471 mental health history.”

472 **Theme VII: Peer Consultation and Training for Clinicians** Two participants discussed the po-
473 tential utility of the model for peer consultation and clinical training. Clinicians began to envision
474 how they might incorporate the model into their decision-making processes, suggesting a growing
475 willingness to see themselves as potential users.

476 Clinician Xing mentioned that if the model’s output differed from their judgment, they would seek
477 guidance from peers or supervisors, acknowledging the importance of human expertise in complex
478 clinical cases:

479 “I am worried about the weight of some questions which are not the signature of depression,
480 but it also adds in personal bias to the tool. I would seek peers/supervisors for a second
481 opinion when the model differs from my clinical judgement.”

482 Clinician Lina discussed that effective training on model tuning is crucial to avoid bias and ensure
483 the model’s appropriate use. Ongoing guidance would help clinicians incorporate evidence-based
484 practices into their work:

486 “Will there be more guidance on using this tool to enhance clinical judgment or incorporate
487 evidence-based criteria? Counselors’ biases can be identified by comparing similar clients’
488 depression symptoms and analyzing their similarities and differences.”

489
490 **Theme VIII: Varied Potentials in the Future** When considering the future of the model, partic-
491 ipants expressed a mix of concerns and excitement. Two participants discussed the importance of
492 having an user-friendly and visual interface to enhance the user experience, allowing clinicians to
493 see the impact of weight changes clearly. This would improve engagement and the overall utility of
494 the model:

495 “I would like to make this model more user friendly, visual, and show more clear
496 changes/interpretation after each feature weights change.”

497 “Ensuring that the tool remains user-friendly and offers high fidelity in its outputs is crucial
498 for its future success. I look forward to seeing how it develops further and hope that with
499 some improvements, it can greatly assist clinicians.”

500
501 Clinician Wang raised concerns about allowing all users to adjust the feature weights, questioning
502 whether this flexibility might compromise accuracy. Other participants echoed this concern, wor-
503 rying that over-manipulation by clinicians could cause important peripheral data to be overlooked,
504 potentially leading to inaccurate diagnoses:

505 “I am not sure what’s the usage of this model. It’s like allowing every clinician to change
506 the algorithm, than how is that applicable to large amount of clients?”

507 “Over-manipulation by the clinician that may overlook the importance of some peripheral
508 data. It’s adjustable but not sure how accurate the result is.”

509
510 We observed differences between doctorate- and master-level clinicians during the interviews.
511 Doctorate-level clinicians were more proactive, asking clarifying questions, adjusting features cre-
512 atively, and exploring how to integrate the model into their work. In contrast, master-level clinicians
513 were more defensive, raising concerns that led to shorter, less in-depth interviews. These differences
514 may be due to doctorate clinicians’ broader roles and theoretical training, while master-level clini-
515 cians, burdened by heavy caseloads and limited support, focused more on practical applications and
516 found it difficult to engage with the prototype without a polished user interface.

517 518 6 CONCLUSION

519
520 We propose NN-kNN for mental disorder diagnosis not as a solution to the inherent challenges of
521 diagnosis, nor solely for its predictive accuracy. All models, including NN-kNN, are biased because
522 they are at most as effective as the data they are trained on. What sets NN-kNN apart is its full
523 interpretability and adjustability, enabling practitioners to detect and correct biases within both the
524 model and the data. By offering the ability to adjust model parameters based on expert knowledge,
525 NN-kNN empowers clinicians to make more informed and ethical decisions. Practitioners can man-
526 ually adjust the model parameters, discover unforeseen patterns or biases, and decide when it is
527 appropriate to rely on AI in clinical settings.

528 Through our qualitative interviews, practitioners expressed appreciation for the model’s trans-
529 parency, flexibility, and the way it integrates their clinical expertise into the diagnostic process.
530 This balance of AI-driven insights and human judgment has the potential to build greater trust and
531 utility in AI-based diagnostic tools. Looking ahead, we plan to extend NN-kNN by incorporating
532 complex data from additional modalities, further supporting clinicians in delivering more holistic
533 and data-rich diagnoses.

534 AUTHOR CONTRIBUTIONS

535
536 Hidden for Anonymity

537 538 ACKNOWLEDGMENTS

539
540 Hidden for Anonymity

REFERENCES

- 540
541
542 Juan Aguilera, Luis C. González, Manuel Montes-y Gómez, and Paolo Rosso. A new weighted
543 k-nearest neighbor algorithm based on newton’s gravitational force. In Ruben Vera-Rodriguez,
544 Julian Fierrez, and Aythami Morales (eds.), *Progress in Pattern Recognition, Image Analysis,
545 Computer Vision, and Applications*, pp. 305–313, Cham, 2019. Springer International Publishing,
546 ISBN 978-3-030-13469-3.
- 547 David W. Aha and Robert L. Goldstone. Concept learning and flexible weighting. In *In Proceedings
548 of the Fourteenth Annual Conference of the Cognitive Science Society*, pp. 534–539. Erlbaum,
549 1992.
- 550 American Psychiatric Association. *Diagnostic and statistical manual of mental disorders: DSM-5*,
551 volume 5. American psychiatric association Washington, DC, 2013.
- 552
553 Sheena K Au-Yeung, Louise Bradley, Ashley E Robertson, Rebecca Shaw, Simon Baron-Cohen,
554 and Sarah Cassidy. Experience of mental health diagnosis and perceived misdiagnosis in autistic,
555 possibly autistic and non-autistic adults. *Autism*, 23(6):1508–1518, 2019.
- 556
557 Getinet Ayano, Sileshi Demelash, Zegeye Yohannes, Kibrom Haile, Mikiyas Tulu, Dawit Assefa,
558 Abel Tesfaye, Kelemua Haile, Melat Solomon, Asrat Chaka, et al. Misdiagnosis, detection rate,
559 and associated factors of severe psychiatric disorders in specialized psychiatry centers in ethiopia.
560 *Annals of general psychiatry*, 20:1–10, 2021.
- 561
562 Gautam Bhattacharya, Koushik Ghosh, and Ananda S. Chowdhury. Granger causality driven ahp
563 for feature weighted knn. *Pattern Recogn.*, 66(C):425–436, June 2017. ISSN 0031-3203. doi: 10.
564 1016/j.patcog.2017.01.018. URL <https://doi.org/10.1016/j.patcog.2017.01.018>.
- 565
566 M. Bicego and M. Loog. Weighted k-nearest neighbor revisited. In *2016 23rd International Confer-
567 ence on Pattern Recognition (ICPR)*, pp. 1642–1647, 2016. doi: 10.1109/ICPR.2016.7899872.
- 568
569 Nimagna Biswas, Saurajit Chakraborty, Sankha Subhra Mullick, and Swagatam Das. A parameter
570 independent fuzzy weighted k-nearest neighbor classifier. *Pattern Recognition Letters*, 101:80–
571 87, 2018.
- 572
573 Andrea Bonzano, Pádraig Cunningham, and Barry Smyth. Using introspective learning to improve
574 retrieval in CBR: A case study in air traffic control. In D. Leake and E. Plaza (eds.), *Proceedings of
575 the 2nd International Conference on Case-Based Reasoning (ICCBR-97)*, volume 1266 of *LNAI*,
576 pp. 291–302, Berlin, July 25–27 1997. Springer. ISBN 3-540-63233-6.
- 577
578 Danilo Bzdok and Andreas Meyer-Lindenberg. Machine learning for precision psychiatry: Oppor-
579 tunities and challenges. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, 3(3):
580 223–230, 2018. ISSN 2451-9022. doi: <https://doi.org/10.1016/j.bpsc.2017.11.007>. URL <https://www.sciencedirect.com/science/article/pii/S2451902217302069>.
- 581
582 Subhagata Chattopadhyay. A neuro-fuzzy approach for the diagnosis of depression. *Applied Com-
583 puting and Informatics*, 13(1):10–18, 2017. ISSN 2210-8327. doi: <https://doi.org/10.1016/j.aci.2014.01.001>. URL <https://www.sciencedirect.com/science/article/pii/S2210832714000027>.
- 584
585 Arun Das and Paul Rad. Opportunities and challenges in explainable artificial intelligence (xai): A
586 survey. *arXiv preprint arXiv:2006.11371*, 2020.
- 587
588 Belur. V. Dasarathy. Nearest neighbor (nn) norms : Nn pattern classification tech-
589 niques. *IEEE Computer Society Tutorial*, 1991. URL <https://cir.nii.ac.jp/crid/1572261550010307072>.
- 590
591 Virginia Eatough and Jonathan A Smith. Interpretative phenomenological analysis. *The Sage hand-
592 book of qualitative research in psychology*, pp. 193–209, 2017.
- 593
Jerome H. Friedman. Flexible metric nearest neighbor classification. Technical report, Stanford
University, 1994.

- 594 Pedro J. García-Laencina, José-Luis Sancho-Gómez, Aníbal R. Figueiras-Vidal, and Michel Verley-
595 sen. K nearest neighbours with mutual information for simultaneous classification and missing
596 data imputation. *Neurocomputing*, 72(7):1483–1493, 2009. ISSN 0925-2312. doi: [https://doi.org/](https://doi.org/10.1016/j.neucom.2008.11.026)
597 [10.1016/j.neucom.2008.11.026](https://doi.org/10.1016/j.neucom.2008.11.026). URL [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S0925231209000149)
598 [article/pii/S0925231209000149](https://www.sciencedirect.com/science/article/pii/S0925231209000149). *Advances in Machine Learning and Computational*
599 *Intelligence*.
- 600 Lawrence Gates and David B. Leake. Evaluating cbr explanation capabilities: Survey and next steps.
601 In *ICCBR Workshops*, 2021. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:244731832)
602 [244731832](https://api.semanticscholar.org/CorpusID:244731832).
- 603 Jacob Goldberger, Geoffrey E Hinton, Sam Roweis, and Russ R Salakhutdinov. Neighbour-
604 hood components analysis. In *Advances in Neural Information Processing Systems*, volume 17.
605 MIT Press, 2004. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2004/file/42fe880812925e520249e808937738d2-Paper.pdf)
606 [2004/file/42fe880812925e520249e808937738d2-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2004/file/42fe880812925e520249e808937738d2-Paper.pdf).
- 607 Sarah Graham, Colin Depp, Ellen E Lee, Camille Nebeker, Xin Tu, Ho-Cheol Kim, and Dilip V
608 Jeste. Artificial intelligence for mental health and mental illnesses: an overview. *Current psychi-*
609 *atry reports*, 21:1–18, 2019.
- 610 Ngumimi Karen Iyortsuun, Soo-Hyung Kim, Min Jhon, Hyung-Jeong Yang, and Sudarshan Pant. A
611 review of machine learning and deep learning approaches on mental health diagnosis. In *Health-*
612 *care*, volume 11, pp. 285. MDPI, 2023.
- 613 Hanna Jarvie and Hanna Lindén. Exploring human therapists’ perspectives on artificial intelligence
614 therapists in mental health care, 2024.
- 615 Dan W Joyce, Andrey Kormilitzin, Katharine A Smith, and Andrea Cipriani. Explainable artificial
616 intelligence for mental health through transparency and interpretability for understandability. *npj*
617 *Digital Medicine*, 6(1):6, 2023.
- 618 Elma Kerz, Sourabh Zanwar, Yu Qiao, and Daniel Wiechmann. Toward explainable ai (xai) for
619 mental health detection based on language behavior. *Frontiers in psychiatry*, 14:1219479, 2023.
- 620 Clinton Lau, Xiaodan Zhu, and Wai-Yip Chan. Automatic depression severity assessment with deep
621 learning using parameter-efficient tuning. *Frontiers in Psychiatry*, 14:1160291, 2023.
- 622 Elizabeth Li, David Kealy, Katie Aafjes-van Doorn, James McCollum, John T Curtis, Xiaochen Luo,
623 and George Silberschatz. “it felt like i was being tailored to the treatment rather than the treatment
624 being tailored to me”: Patient experiences of helpful and unhelpful psychotherapy. *Psychotherapy*
625 *Research*, pp. 1–15, 2024.
- 626 Oscar Li, Hao Liu, Chaofan Chen, and Cynthia Rudin. Deep learning for case-based reasoning
627 through prototypes: A neural network that explains its predictions. In *Proceedings of the Thirty-*
628 *Second AAAI Conference on Artificial Intelligence*, pp. 3530–3537. AAAI Press, 2018. ISBN
629 978-1-57735-800-8.
- 630 Huabing Liu, Y Joel Wong, Nancy Goodrich Mitts, PF Jonah Li, and Jacks Cheng. A phenom-
631 enological study of east asian international students’ experience of counseling. *International Journal*
632 *for the Advancement of Counselling*, 42:269–291, 2020.
- 633 Youness Manzali, Khalidou Abdoulaye Barry, Rachid Flouchi, Youssef Balouki, and Mohamed
634 Elfar. A feature weighted k-nearest neighbor algorithm based on association rules. *Journal of*
635 *Ambient Intelligence and Humanized Computing*, pp. 1–14, 2024.
- 636 Elena Marchiori. *Class Dependent Feature Weighting and K-Nearest Neighbor Classifica-*
637 *tion*, pp. 69–78. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013. ISBN 978-3-642-
638 39159-0. doi: [10.1007/978-3-642-39159-0\](https://doi.org/10.1007/978-3-642-39159-0_7)
639 [7](https://doi.org/10.1007/978-3-642-39159-0_7). URL [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-642-39159-0_7)
640 [978-3-642-39159-0_7](https://doi.org/10.1007/978-3-642-39159-0_7).
- 641 Joseph A Maxwell. Why qualitative methods are necessary for generalization. *Qualitative Psychol-*
642 *ogy*, 8(1):111, 2021.

- 648 Joni L Mihura, Manali Roy, and Robert A Graceffo. Psychological assessment training in clinical
649 psychology doctoral programs. *Journal of Personality assessment*, 99(2):153–164, 2017.
650
- 651 Christoph Molnar, Giuseppe Casalicchio, and Bernd Bischl. Interpretable machine learning – a brief
652 history, state-of-the-art and challenges. In Irena Koprinska, Michael Kamp, Annalisa Appice,
653 Corrado Loglisci, Luiza Antonie, Albrecht Zimmermann, Riccardo Guidotti, Özlem Özgöbek,
654 Rita P. Ribeiro, Ricard Gavaldà, João Gama, Linara Adilova, Yamuna Krishnamurthy, Pedro M.
655 Ferreira, Donato Malerba, Ibéria Medeiros, Michelangelo Ceci, Giuseppe Manco, Elio Masciari,
656 Zbigniew W. Ras, Peter Christen, Eirini Ntoutsi, Erich Schubert, Arthur Zimek, Anna Monreale,
657 Przemyslaw Biecek, Salvatore Rinzivillo, Benjamin Kille, Andreas Lommatzsch, and Jon Atle
658 Gulla (eds.), *ECML PKDD 2020 Workshops*, pp. 417–431, Cham, 2020. Springer International
659 Publishing. ISBN 978-3-030-65965-3.
- 660 Mauricio Gabriel Orozco-del Castillo, Esperanza Carolina Orozco-del Castillo, Esteban Brito-
661 Borges, Carlos Bermejo-Sabbagh, and Nora Cuevas-Cuevas. An artificial neural network for
662 depression screening and questionnaire refinement in undergraduate students. In Miguel Félix
663 Mata-Rivera and Roberto Zagal-Flores (eds.), *Telematics and Computing*, pp. 1–13, Cham, 2021.
664 Springer International Publishing. ISBN 978-3-030-89586-0.
- 665 Amorette Perkins, Joseph Ridler, Daniel Browes, Guy Peryer, Caitlin Notley, and Corinna Hack-
666 mann. Experiencing mental health diagnosis: a systematic review of service user, clinician, and
667 carer perspectives across clinical settings. *The Lancet Psychiatry*, 5(9):747–764, 2018.
668
- 669 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ” why should i trust you?” explaining the
670 predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference*
671 *on knowledge discovery and data mining*, pp. 1135–1144, 2016.
672
- 673 Francesco Ricci and Paolo Avesani. *Learning a local similarity metric for case-based reasoning*, pp.
674 301–312. Springer Berlin Heidelberg, Berlin, Heidelberg, 1995. ISBN 978-3-540-48446-2. doi:
675 10.1007/3-540-60598-3\27. URL https://doi.org/10.1007/3-540-60598-3_27.
- 676 Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and
677 use interpretable models instead. *Nature Machine Intelligence*, 1:206–215, 05 2019. doi: 10.
678 1038/s42256-019-0048-x.
- 679 Waddah Saeed and Christian Omlin. Explainable ai (xai): A systematic meta-survey of current chal-
680 lenges and future opportunities. *Knowledge-Based Systems*, 263:110273, 2023. ISSN 0950-7051.
681 doi: <https://doi.org/10.1016/j.knosys.2023.110273>. URL <https://www.sciencedirect.com/science/article/pii/S0950705123000230>.
- 682
683
- 684 Rabia Saleem, Bo Yuan, Fatih Kurugollu, Ashiq Anjum, and Lu Liu. Explaining deep neural
685 networks: A survey on the global interpretation methods. *Neurocomputing*, 513:165–180,
686 2022. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2022.09.129>. URL <https://www.sciencedirect.com/science/article/pii/S0925231222012218>.
687
- 688 Jakob M Schoenborn, Rosina O Weber, David W Aha, Jorg Cassens, and Klaus-Dieter Althoff.
689 Explainable case-based reasoning: a survey. In *AAAI-21 Workshop Proceedings*, 2021.
690
- 691 Lloyd S Shapley. A value for n-person games. *Contribution to the Theory of Games*, 2, 1953.
692
- 693 Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Ad-
694 vances in neural information processing systems*, 30, 2017.
- 695 Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one
696 shot learning. *Advances in neural information processing systems*, 29, 2016.
697
- 698 Kilian Q. Weinberger and Lawrence K. Saul. Distance metric learning for large margin nearest
699 neighbor classification. *J. Mach. Learn. Res.*, 10:207–244, jun 2009. ISSN 1532-4435.
700
- 701 Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. *CoRR*, abs/1410.3916, 2014.
URL <https://api.semanticscholar.org/CorpusID:2926851>.

702 D. Wettschereck, D. Aha, and T. Mohri. A review and empirical evaluation of feature-weighting
703 methods for a class of lazy learning algorithms. *Artificial Intelligence Review*, 11(1-5):273–314,
704 February 1997.

705
706 Zachary Wilkerson, David Leake, and David Crandall. Leveraging shap and cbr for dimensionality
707 reduction on the psychology prediction dataset. In *ICCBR Workshops*, pp. 236–240, 2022.

708 Xiaomeng Ye, David Leake, Yu Wang, Ziwei Zhao, and David Crandall. Towards network im-
709 plementation of cbr: Case study of a neural network k-nn algorithm. In *Case-Based Reasoning*
710 *Research and Development: 32nd International Conference, ICCBR 2024, Merida, Mexico, July*
711 *1–4, 2024, Proceedings*, pp. 354–370, Berlin, Heidelberg, 2024. Springer-Verlag. ISBN 978-3-
712 031-63645-5. doi: 10.1007/978-3-031-63646-2_23. URL [https://doi.org/10.1007/
713 978-3-031-63646-2_23](https://doi.org/10.1007/978-3-031-63646-2_23).

714 Tianlin Zhang, Annika M Schoene, Shaoxiong Ji, and Sophia Ananiadou. Natural language pro-
715 cessing applied to mental illness detection: a narrative review. *NPJ digital medicine*, 5(1):1–13,
716 2022.

717 718 719 A APPENDIX

720
721 You may include other additional sections here.

722
723 Authors may use as many pages of appendices (after the bibliography) as they wish, but reviewers
724 are not required to read the appendix.

725 726 B META-PARAMETER SETTINGS

727
728 We performed a standard 10-fold cross-validation on the dataset, training the model until test set
729 accuracy plateaued for a fixed number of epochs (40). Due to the small dataset size and the fact
730 that the goal of this experiment is not to precisely gauge the model’s prediction accuracy, we did not
731 allocate a separate validation set. The model was trained in batches of 4 using the Adam optimizer
732 with a learning rate of 0.01. The top- k case selection layer was disabled, as enabling it led to lower
733 accuracy. Although some parameter choices (e.g., the relatively high learning rate) may seem un-
734 conventional in typical machine learning setups, they are justified and discussed in Ye et al. (2024).
735 Most importantly, additional configurations were employed to enhance the model’s interpretability
736 for practitioners:

- 737 • All cases share the same feature weights ($w_{x\delta_i} = w_{y\delta_i} = w_{\delta_i}$ for any two cases x and y)
738 to reduce the number of parameters. This approach applies global feature weighting, where
739 each feature is assigned a single weight across all cases.
- 740 • For each case x , we initialize the feature weights as $w_{x\delta_i} = 1$, the case bias as $b_x = 50$, the
741 case weight as $w_{(x,L)} = 1$, and all class biases as $b_L = 1$.

742 743 744 C DATASET DESCRIPTION

745
746 The dataset is a collection of answers from 157 undergraduate students to a survey of true/false
747 questions. The questions are designed according to the Diagnostic and Statistical Manual of Mental
748 Disorders by a psychologist (Orozco-del Castillo et al., 2021). The questions are listed below.

- 749 1. Most of the time I have difficulty concentrating on simple tasks
- 750 2. I don’t feel like doing my daily duties
- 751 3. My friends or family have told me that I look different
- 752 4. When I think about the future it is difficult for me to imagine it clearly
- 753 5. People around me often ask me how I feel
- 754 6. I consider that my life is full of good things
- 755

- 756 7. My hobbies are still important to me
757
758 8. I'm still as punctual as I have always been
759 9. If I had the chance, I would spend all day in my bed
760 10. I have found that I can spend a lot of time scrolling the screen of my cell phone without
761 searching or stopping at anything in particular
762 11. When someone asks me something, I have noticed that I take longer than normal to respond
763
764 12. I have noticed my body shaken without any cause
765 13. I felt more encouraged to do my daily activities before
766 14. Sometimes I wake up sad and I can't explain why
767 15. In recent months I usually reproach myself for things from the past
768 16. I think my thoughts are strange or different from before
769 17. I feel guilty about the decisions that I have made
770
771 18. I don't feel as comfortable with my body as I did before
772 19. I don't feel successful compared to others
773 20. It is difficult for me to make decisions even if they are simple
774 21. I'm capable of achieving what I propose to myself
775
776 22. It is not difficult for me to understand something the first time
777 23. I have thought more than before about what my death would be like
778 24. Being dead seems to be a solution to some problems
779 25. I would rather stay home than go out with my friends
780
781 26. I like to attend family gatherings
782 27. I feel excited when thinking about my life project
783 28. The decisions I have made so far have been the right ones
784 29. I am able to carry out my activities as I have always been
785 30. I like to be in touch with my friends and family through social media
786
787 31. It is easy for me to choose a photograph of myself to show it on social media
788 32. I am proud of what I have achieved so far
789 33. I have trouble remembering things easily
790
791 34. In recent months I have had discussions with my schoolmates or colleagues
792 35. I constantly imagine that something will go wrong at my work or at school
793 36. I am afraid of being wrong when doing my homework
794 37. I'm not too worried about what might happen in a few weeks
795 38. Lately it's hard for me to calm down
796 39. Everything will be alright
797 40. I can easily blank my mind
798
799 41. I am bothered by insignificant things that were not important before
800 42. I find it uncomfortable to be in a crowded place
801 43. Sometimes I feel trapped
802 44. I am easily frightened by unexpected noises
803 45. I have difficulties to do one task at a time
804 46. I have the feeling that I am forgetting to do something
805 47. I can clearly express to others how I feel
806 48. I can sleep easily
807 49. I enjoy every moment of the day
808
809

- 810 50. I imagine that at any moment a disaster of nature may occur
811
812 51. Sometimes I feel like I get tired easily
813
814 52. Being locked in an elevator would be the worst thing that could happen to me
815
816 53. I'm bothered by people walking slowly in front of me
817
818 54. I don't usually get upset if something doesn't go as expected
819
820 55. Sometimes it is as if some conversations with friends or family become interrogations
821
822 56. I manage my schedule as I always have
823
824 57. It bothers me to feel that people on the street approach me
825
826 58. I have no difficulty understanding what people explain to me
827
828 59. I consider that I am good at controlling my emotions
829
830 60. In new situations I feel calm and encouraged
831
832 61. Sometimes I forget what I wanted to say because I have several thoughts at the same time
833
834 62. I would like to know what will happen in the future
835
836 63. When I get angry I can easily explode
837
838 64. I can put down my cell phone and dedicate myself to reading without distractions
839
840 65. I worry that people will not understand what I mean
841
842 66. Sometimes I do not listen to what people say to me because I am thinking about other
843 things
844
845 67. I get angry easily
846
847 68. I'm afraid that something bad could happen to me
848
849 69. It is not important for me to meet set dates
850
851 70. I like to think clearly before giving my opinion
852
853 71. I use lies just to get out of certain problems
854
855 72. If I have the opportunity to get in line to avoid wasting time, I do it
856
857 73. I have difficulty making elaborate plans
858
859 74. People have problems because of themselves
860
861 75. No me parece importante lo que los otros piensen sobre mí
862
863 76. Laws are not as important as others think
77. I would regret betraying a friend
78. I prefer that a negotiation supports the largest possible number of people involved
79. It is easy for me to work in a team
80. It is important to help people when they need it
81. I have punched someone or thought of doing it
82. If it was necessary I would pretend to be someone else to get something
83. I consider it important to ensure my physical safety and that of those around me
84. After an argument I usually go over what happened in my head
85. I have a hard time controlling myself when I get angry
86. Loyalty is important
87. If I can help a person I will stop what I'm doing to help them
88. Sometimes people need physical force to understand
89. It makes me laugh when my superiors at school or at work demand something
90. Deceiving people is not wrong if it is to achieve something important
91. I like to greet my neighbors
92. It does not seem serious to me to have some debts

- 864 93. People steal because they have needs
865 94. I lose control easily
866 95. Neighbors must put up with each other’s noises without complaining
867 96. Littering on public roads is wrong
868 97. People who commit crimes have their reasons for doing it
869 98. It is normal to change jobs several times a year
870 99. It is important to respect turns
871 100. I could pretend to be someone else to achieve what I want
872 101. I consider it important that all people have the same rights
873 102. I have a hard time taking ”no” for an answer
874
875
876

877 D DATA PREPROCESSING

879 The dataset was used to train neural networks for depression screening and then later used for an
880 explainable AI challenge in the Explainable AI Challenge at the 2022 International Conference on
881 Case-Based Reasoning (Wilkerson et al., 2022). We obtained a version of the dataset with 104
882 cases with 102 attributes and a risk score (from 1 to 5) calculated from the number of physical
883 symptoms related to depression. We oversampled less frequent classes to counteract class imbalance
884 and transformed risk scores into three classes (low, medium, and high risk of depression) following
885 the example of Wilkerson et al. (2022). Our final dataset contains 117 cases; Each case has 102
886 features of binary values and a class label of 0, 1 or 2, representing low, medium and high risk.
887

888 E INTERVIEW QUESTIONS

889 Following are the questions we used when interviewing the 10 participating practitioners.
890
891

892 E.1 DEMOGRAPHIC QUESTIONNAIRE

- 893 1. What’s your full name?
894 2. What’s your gender?
895 3. What’s your age?
896 4. What’s your race?
897 5. What’s your sexual orientation?
898 6. What’s your highest degree and for how long have you been licensed?
899 7. What’s your professional field (Counseling, clinical, social work, school psychology, etc.)?
900 8. What’s your current job?
901 9. Which state are you currently living in?
902 10. What are your primary clinical populations?
903 11. What are your primary theoretical orientations?
904 12. Please rate your knowledge about AI technology from 0-10 (0 being absolutely no knowl-
905 edge, 10 being complete expertise).
906
907
908
909

910 E.2 QUALTRICS QUALITATIVE QUESTIONS

911 **As the preliminary user:**

- 912
913 1. Regarding the tunable experience you just had, did you feel like the model has become
914 more useful clinically after you tuned the feature weights? Please elaborate.
915 2. Did you feel like the model has become more clinically trustworthy? Please elaborate.
916 3. Please provide any other comments or thoughts about your experience trying out our algo-
917 rithm just now.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

As a practitioner thinking about using this model for clinical diagnosis:

1. If you were to use this AI algorithm, would the tunable feature be useful for you?
2. What would be the strengths and concerns you have about using this tunable feature?
3. Since our model can detect bias, if the model’s explanations differ from your clinical judgment, what would you do?
4. Since our model is fully interpretable and tunable, would this algorithm be more ethically trustworthy?
5. What other factors should we consider in improving this model for clinical diagnosis?