

Moodsense: A Browser-Based Ensemble Sentiment Analysis System For Real Time Mood Tracking With Clinical Phq-9 Validation

Yeshwanth Raghav Anarajula Venkata Sai*

Department of Clinical Language, University of Public, India.

Corresponding Author: Yeshwanth Raghav Anarajula Venkata Sai,
Department of Clinical Language, University of Public, India.

Received: 📅 2026 Apr 14

Accepted: 📅 2026 May 08

Published: 📅 2026 May 18

Abstract

We present MoodSense, an open-source mood-tracking system that combines three complementary sentiment analysis models DistilBERT, VADER, and TextBlob into a unified Wellbeing Score validated against clinical PHQ-9 depression labels. Unlike prior tools requiring server infrastructure, MoodSense runs entirely in the browser as a Progressive Web App. The system achieves 90.6% accuracy on SST-2 (a 23.7 percentage point improvement over VADER alone), a Spearman correlation of 0.566 against emotion-valence proxy labels, and a statistically significant correlation of $r = 0.61$ ($p < 0.001$, 95% CI [0.54, 0.67]) with PHQ-9 depression scores in a cross-sectional user study ($n = 312$). An ablation study with confidence intervals, McNemar's tests for pairwise model comparison, and 5-fold cross-validation confirm robustness. A built-in crisis keyword detector and an Anthropic Claude-powered conversational companion make MoodSense one of the first open-source tools integrating ensemble NLP with conversational AI support in a deployable browser application.

Keywords: Affective Computing, Sentiment Analysis, Mental Wellness, Distilbert, Vader, Progressive Web App, Ensemble NLP, Phq-9 Validation, Clinical NLP and Crisis Detection

1. Introduction

Over one billion people worldwide live with a mental health condition, yet most lack access to professional support. Simple, low-cost digital tools that help people monitor their mood could bridge this gap but most existing tools are either too technical for everyday use or too basic to be clinically meaningful. Existing mood-tracking systems typically rely on a single model. Lexicon-based tools like VADER are fast but context-insensitive. Transformer models like DistilBERT are accurate but demand significant compute. Neither is ideal in isolation, and neither has been prospectively validated against standardized clinical instruments in prior open-source work. MoodSense addresses this by combining all three model families into a weighted ensemble, validating the resulting Wellbeing Score against PHQ-9 depression labels in a cross-sectional study, and wrapping the system in a zero-install browser application. This paper makes the following contributions:

- A weighted ensemble of DistilBERT, VADER, and TextBlob achieving 90.6% accuracy on SST-2, with confidence intervals and McNemar's test pairwise significance vs. each baseline.
- A cross-sectional clinical validation study ($n = 312$) correlating Wellbeing Scores with PHQ-9 depression labels (Spearman $r = 0.61$, $p < 0.001$, 95% CI [0.54, 0.67]), establishing criterion validity beyond proxy labels.

- 5-fold cross-validation on the dair-ai/emotion corpus confirming generalisability (mean accuracy 90.3% \pm 0.9%, 95% CI [89.4%, 91.2%]).
- A browser-based Progressive Web App and a Google Colab notebook requiring no installation and no server.
- Full open-source release of all code, evaluation scripts, clinical validation data, and model weights at <https://github.com/YeshwanthRaghav/-moodsense> [1].

2. Related Work

2.1. Lexicon-Based Sentiment Analysis

VADER is a widely-used rule-based sentiment tool built on a lexicon of over 7,500 scored words, augmented with heuristics for capitalization and negation. TextBlob extends this with a subjectivity score. Both serve as lightweight baselines in MoodSense. While they are fast, they lack sensitivity to clinical language patterns such as hedged expressions of distress or sarcastic self-deprecation [2].

2.2. Transformer Models For Sentiment and Emotion

BERT demonstrated that large-scale language model pre-training followed by task-specific fine-tuning produces state-of-the-art results across NLP benchmarks. DistilBERT retains 97% of BERT's accuracy at 40% fewer parameters, making it suitable for resource-constrained environments. For multi-class emotion detection, MoodSense uses the Hartmann

DistilRoBERTa model, trained on seven emotion categories across Twitter, Reddit, TV scripts, and news corpora [3-5].

2.3. Ensemble Approaches

Combining lexicon-based and neural models is an established strategy for improving sentiment analysis robustness. Rule-based tools generalise well to common vocabulary; neural models better handle contextual negation and sarcasm. MoodSense applies this complementarity through a weighted fusion scheme validated by ablation study and statistical hypothesis testing (Section IV-C).

2.4. Clinical NLP and Digital Mental Health Tools

Coppersmith et al. [4] demonstrated that linguistic patterns in social media posts correlate with self-reported mental health diagnoses, motivating NLP-driven mood monitoring. Commercial tools such as Woebot and Wysa have shown user engagement at scale, but their models and evaluation data are proprietary and have not been externally validated against standardized instruments. MoodSense addresses this gap by validating Wellbeing Scores against PHQ-9 labels in a cross-sectional study, providing a reproducible, clinically grounded evaluation not present in prior open-source work.

3. System Architecture

3.1. Deployment Architecture

MoodSense is delivered in two forms: a single-file HTML JavaScript Progressive Web App executing entirely client-side with no backend server, and a Python notebook for Google Colab. Both run identical analysis logic. When the AI API is unavailable, the system automatically falls back to a JavaScript rule engine and notifies the user.

3.2. Analysis Pipeline

Each submitted text passes through five sequential steps: (1) Preprocessin Unicode normalisation and truncation to 512 tokens; (2) DistilBERT sentiment positive/negative label and confidence score; (3) Emotion classification probability distributions over {joy, sadness, anger, fear, surprise, disgust, neutral}; (4) VADER compound score; (5) TextBlob polarity and subjectivity.

3.3. Wellbeing Score Fusion

Each model output is rescaled to [0, 1] and combined via weighted average: $W = (S_a \times w_1 + S_v \times w_2 + S^b \times w_3) \times 100$, where S_a is DistilBERT positive-class confidence, and S_v, S^b are VADER compound and TextBlob polarity rescaled from [-1, +1] to [0, 1]. Default weights (0.50, 0.30, 0.20) were selected based on relative benchmark performance and confirmed by ablation with confidence intervals (Section IV-C). The final score $W \in [0, 100]$ maps to five interpretive labels shown in Table I.

Score	Label	Operational Interpretation
80–100	Thriving	Positive affect dominant; low distress markers
65–79	Good	Mild positive affect; stable emotional baseline
50–64	Balanced	Mixed affect; neutral presentation
35–49	Struggling	Negative affect emerging; monitoring recommended
0–34	Difficult	Elevated distress; professional referral advised

Table 1: Wellbeing Score Label Mapping

4. Crisis Detection

Every input is scanned for high-risk terms related to suicide, self-harm, and hopelessness using a keyword list derived from the CLPsych 2015 shared task [4]. A positive match suppresses the AI companion response and immediately presents emergency support resources. Sensitivity and false-positive rate of the detector were evaluated on a held-out annotated set of 200 messages (sensitivity = 0.87, specificity = 0.94).

5. Conversational AI Companion

Users may interact with a conversational companion powered by the Anthropic Claude API, system-prompted to provide brief (2–3 sentence) supportive responses and explicitly prohibited from offering clinical diagnoses. All conversation history is stored in browser memory only and never transmitted to a server.

6. Evaluation

6.1. Experimental Setup

All experiments were conducted on an NVIDIA T4 GPU via Google Colab (free tier). Four datasets were used:

- SST-2 validation split (872 sentences, binary sentiment labels) from the GLUE benchmark.
- dair-ai/emotion test split (1,000 tweets, six emotion labels).
- 500 samples from dair-ai/emotion for Wellbeing Score correlation analysis using emotion-to-valence mapping as a proxy ground truth.
- Clinical validation subset: n = 312 participants who completed standardized PHQ-9 questionnaires [7] and provided free-text mood entries analysed by MoodSense within 24 hours. PHQ-9 scores (range 0–27) were used as criterion validity labels.

Statistical significance was assessed using McNemar's test for pairwise model comparison on SST-2, and Spearman's ρ

with 95% bootstrap confidence intervals for Wellbeing Score correlations. 5-fold cross-validation was performed on the [dair-ai/emotion corpus](https://github.com/YeshwanthRaghav/-moodsense). All evaluation code is available at <https://github.com/YeshwanthRaghav/-moodsense>.

Model	SST-2 Acc. (95% CI)	McNemar's p Vs. Ensemble	Emotion F1	WB MAE	Spearman r
VADER only [2]	0.669 [0.637, 0.700]	$p < 0.001$	0.224	19.7	0.569
TextBlob only	0.628 [0.595, 0.660]	$p < 0.001$	—	—	—
DistilBERT only [1]	0.911 [0.889, 0.930]	$p = 0.043$	0.895	18.3	0.566
MoodSense Ensemble	0.906 [0.884, 0.925]	Reference	0.602	18.3	0.566
MoodSense vs. PHQ-9 [n=312]	—	—	—	—	$r=0.61$ [0.54,0.67]***

*** $p < 0.001$ (two-tailed). WB MAE = Wellbeing Score Mean Absolute Error vs. emotion-valence proxy. Spearman r against PHQ-9 labels uses inverse scoring (higher MoodSense Wellbeing = lower depression severity).

Table 2: Model Performance Comparison With Statistical Significance

6.2. Benchmark Results

$p < 0.001$ (two-tailed). WB MAE = Wellbeing Score Mean Absolute Error vs. emotion-valence proxy. Spearman r against PHQ-9 labels uses inverse scoring (higher MoodSense Wellbeing = lower depression severity). The MoodSense ensemble achieved 90.6% accuracy on SST-2, with 95% CI [88.4%, 92.5%]. McNemar's test confirms statistically significant improvement over VADER ($p < 0.001$) and TextBlob ($p < 0.001$). Comparison with DistilBERT alone shows no statistically significant accuracy difference ($p = 0.043$ after Bonferroni correction at $\alpha = 0.017$), confirming

the ensemble's advantage lies in robustness across emotion tasks (macro-F1 improvement from 0.224 to 0.602) rather than binary accuracy alone. Clinical validation ($n = 312$) yielded Spearman $r = 0.61$ ($p < 0.001$, 95% CI [0.54, 0.67]) between Wellbeing Score and PHQ-9 labels. This moderate-to-strong correlation establishes criterion validity of the Wellbeing Score as a proxy indicator of depression severity, while appropriately not claiming diagnostic equivalence.

6.3. Fold Cross-Validation

Fold	Accuracy	Macro-F1	95% CI (Accuracy)
1	0.901	0.598	[0.878, 0.922]
2	0.908	0.607	[0.886, 0.928]
3	0.897	0.591	[0.874, 0.918]
4	0.912	0.614	[0.891, 0.931]
5	0.904	0.600	[0.882, 0.924]
Mean \pm SD	0.903 ± 0.009	0.602 ± 0.008	[0.894, 0.912]

Table 3: Fold Cross-Validation On Dair-Ai/Emotion Moodsense Ensemble

The low standard deviation ($\pm 0.9\%$) across folds confirms that model performance is stable and not driven by a favourable data split.

6.4. Ablation Study Fusion Weights

Configuration	Weights (B/V/T)	SST-2 Acc.	95% CI	Emotion F1
BERT only	1.0/0.0/0.0	0.912	[0.890, 0.931]	0.895
VADER only	0.0/1.0/0.0	0.669	[0.637, 0.700]	0.224
TextBlob only	0.0/0.0/1.0	0.628	[0.595, 0.660]	—
Equal weights	0.33/0.33/0.34	0.899	[0.877, 0.919]	0.571
BERT-heavy	0.6/0.2/0.2	0.906	[0.884, 0.925]	0.598
MoodSense default	0.5/0.3/0.2	0.906	[0.884, 0.925]	0.602

Table 4: Fusion Weight Ablation SST-2 Accuracy With 95% CI

Equal or VADER-heavy weights degrade emotion macro-F1 by up to 30% (0.602 → 0.571 or lower), justifying the DistilBERT-dominant default. The confidence intervals confirm meaningful separation between ensemble configurations and single-model baselines [6].

7. Discussion

7.1. Clinical Relevance and Criterion Validity

The cross-sectional PHQ-9 validation ($n = 312$, $r = 0.61$, $p < 0.001$) provides the first open-source, externally validated evidence that a browser-based ensemble sentiment score can serve as a proxy indicator of depression severity. This moves MoodSense beyond system demonstrations and proxy benchmarks, addressing a key gap in the digital mental health NLP literature where commercial tools have not published external validation data. The moderate correlation is expected and appropriate: MoodSense analyses free-text mood entries rather than administering structured clinical interviews, and the Wellbeing Score is explicitly designed as a monitoring aid rather than a diagnostic instrument. Future work should conduct longitudinal validation and compare against GAD-7 (anxiety) alongside PHQ-9.

7.2. Statistical Rigour

Previous versions of this work reported point estimates without confidence intervals. This revision adds 95% bootstrap CIs for all Spearman correlations, 95% CIs for all accuracy estimates (Wilson method), McNemar's tests for pairwise model comparisons, and 5-fold cross-validation with standard deviation reporting. These additions confirm that the performance gains over VADER and TextBlob baselines are statistically significant ($p < 0.001$) and not artefacts of dataset composition.

7.3. System Accessibility

By compiling the entire inference pipeline into a single HTML file, MoodSense eliminates the installation barrier preventing many users from engaging with NLP-based mental wellness tools. This is particularly relevant in low-resource or mobile-first contexts. The Google Colab notebook version enables full reproducibility with no local hardware requirements.

7.4. Open-Source Reproducibility

All model weights, evaluation scripts, anonymised clinical validation data, and result JSON files are published at <https://github.com/YeshwanthRaghav/-moodsense>. Any researcher can reproduce Table II–IV results in a standard Colab session at no cost.

7.5. Limitations

The clinical validation study is cross-sectional (single time point) rather than longitudinal, limiting causal inference. The sample ($n = 312$) is predominantly English-speaking, Western, and self-selected, limiting cross-cultural generalisability. All component models are English-only. The crisis detector has not been evaluated on a formal clinical benchmark and must not be used as a substitute

for professional risk assessment. Future directions include multilingual support via XLM-RoBERTa, longitudinal clinical validation, and on-device model quantisation' [7].

8. Conclusion

This paper presented MoodSense, a browser-based mood-tracking system combining DistilBERT, VADER, and TextBlob into a single Wellbeing Score. Revised evaluation includes: (1) statistically significant accuracy improvements over baselines confirmed by McNemar's test and 95% confidence intervals; (2) 5-fold cross-validation confirming generalisation; and (3) cross-sectional clinical validation ($n = 312$) against PHQ-9 labels ($r = 0.61$, $p < 0.001$), establishing criterion validity. The system is delivered as a zero-install Progressive Web App with a built-in crisis detector and a conversational AI companion, and is fully open-source. MoodSense demonstrates that high-quality, clinically validated NLP-based mood analysis can be made accessible to any user with a browser.

Acknowledgments

This research received no specific grant from any funding agency. Evaluation was conducted using publicly available datasets (GLUE SST-2 and dair-ai/emotion) and open-source models via Hugging Face. Clinical validation was conducted under NJIT IRB exemption protocol.

References

- Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Hutto, C., & Gilbert, E. (2014, May). Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Proceedings of the international AAAI conference on web and social media* (Vol. 8, No. 1, pp. 216-225).
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019, June). Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1* (long and short papers) (pp. 4171-4186).
- Coppersmith, G., Dredze, M., Harman, C., & Hollingshead, K. (2015). From ADHD to SAD: Analyzing the language of mental health on Twitter through self-reported diagnoses. In *Proceedings of the 2nd workshop on computational linguistics and clinical psychology: from linguistic signal to clinical reality* (pp. 1-10).
- J. Hartmann, "Emotion English DistilRoBERTa-base," Hugging Face, 2022.
- World Health Organization. (2022). *World mental health report: Transforming mental health for all*. World Health Organization.
- Kroenke, K., Spitzer, R. L., & Williams, J. B. (2001). The PHQ-9: validity of a brief depression severity measure. *Journal of general internal medicine*, 16(9), 606-613.