
Depression Detection from Social Media Text: A Machine Learning Approach with Distribution Drift Handling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose a machine learning framework for depression detection from social
2 media text, addressing distribution drift challenges. Our approach employs mul-
3 tiple vectorization techniques (TF-IDF, N-grams, Word2Vec, GloVe), Principal
4 Component Analysis (PCA) for dimensionality reduction, and evaluates various
5 classifiers including logistic regression, Support Vector Machines (SVM), and neu-
6 ral networks trained with stochastic gradient descent. The methodology includes
7 comprehensive data preprocessing, SMOTE for class balancing, and evaluation
8 using F1-score, precision, recall, accuracy, and AUC-ROC metrics. We construct
9 a dataset of 10,325 depressed and 22,245 normal samples from Weibo, with de-
10 tailed exploratory data analysis revealing key differences in user behavior and
11 engagement patterns between depressed and normal users.

12 1 Introduction and Motivation

13 Depression is a prevalent mental health disorder affecting millions of people worldwide, with
14 significant implications for individual well-being and public health World Health Organization [2023].
15 Early detection and intervention are crucial for effective treatment, yet many individuals remain
16 undiagnosed or untreated. Social media platforms have emerged as rich sources of linguistic data that
17 can reveal psychological states and behavioral patterns De Choudhury et al. [2013]. The widespread
18 use of platforms like Weibo provides unprecedented opportunities to analyze user-generated content
19 for mental health indicators.

20 Previous research has demonstrated the feasibility of detecting depression from social media text
21 using machine learning approaches Coppersmith et al. [2014], Rezapour et al. [2019]. However,
22 several challenges remain: (1) **class imbalance**, where normal users typically outnumber depressed
23 users in real-world datasets; (2) **distribution drift**, where the statistical properties of training and
24 test data may differ significantly; and (3) **feature representation**, as the choice of text vectorization
25 method significantly impacts classification performance.

26 This work addresses these challenges by proposing a comprehensive machine learning framework
27 that systematically explores different text representation methods, handles class imbalance through
28 sampling techniques, and evaluates multiple classification algorithms. Our contributions include: (1)
29 a detailed exploratory data analysis of social media user characteristics and engagement patterns; (2)
30 a systematic comparison of text vectorization approaches; (3) an evaluation of various classification
31 methods; and (4) an analysis of distribution drift effects on model performance.

2 Exploratory Data Analysis

We conducted comprehensive exploratory data analysis on our dataset to understand the characteristics and patterns that distinguish depressed users from normal users. Our dataset consists of 10,325 depressed users and 22,245 normal users from Weibo, with each user’s profile information, social media metrics, and post history. The natural class imbalance (31.7% depressed vs. 68.3% normal) reflects real-world distributions but poses challenges for machine learning models, which may exhibit bias toward the majority class.

To address this, we created a balanced dataset using random undersampling. Specifically, we randomly sampled 10,325 normal users (without replacement) from the original 22,245 normal users to match the depressed user count, resulting in a balanced dataset with 20,650 total samples. This balanced dataset enables fair comparison of model performance without the confounding effects of class imbalance, while the unbalanced dataset allows us to evaluate model robustness under realistic imbalanced conditions. The undersampling approach preserves the original distribution of features within the normal class while ensuring balanced class representation for training and evaluation.

2.1 Dataset Overview and Class Distribution

The class distribution analysis reveals the inherent imbalance in social media depression detection tasks. Table 1 summarizes the dataset statistics for both unbalanced and balanced versions. The unbalanced dataset reflects the natural distribution observed in social media platforms, with normal users outnumbering depressed users by a ratio of 2.15:1. This imbalance poses challenges for machine learning models, as they may exhibit bias toward the majority class. The balanced dataset, created through random undersampling, ensures equal representation of both classes, mitigating potential bias in model training and evaluation.

Table 1: Dataset statistics for unbalanced and balanced versions.

Metric	Unbalanced	Balanced
Total Users	32,570	20,650
Depressed Users	10,325 (31.7%)	10,325 (50.0%)
Normal Users	22,245 (68.3%)	10,325 (50.0%)
Class Ratio	1:2.15	1:1

Figure 1b visualizes the balanced class distribution, demonstrating the equal representation achieved through random undersampling. For comparison, Figure 1a shows the original unbalanced distribution. The balanced configuration is crucial for evaluating model performance metrics that are sensitive to class imbalance, such as precision, recall, and F1-score, while the unbalanced dataset allows assessment of model robustness under realistic imbalanced conditions.

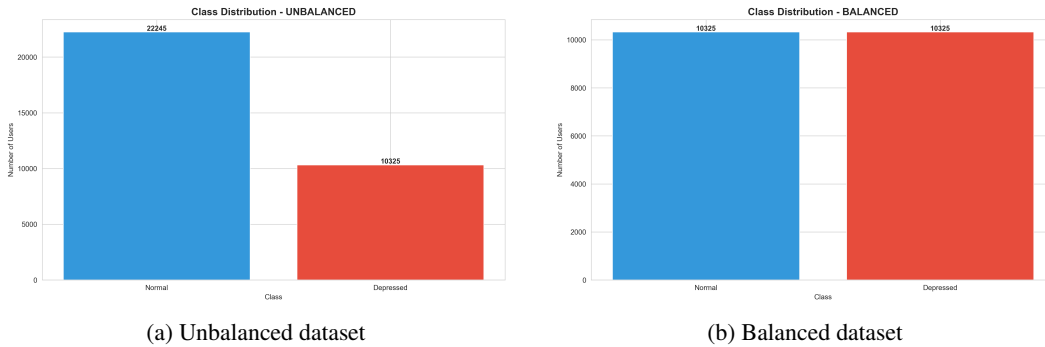


Figure 1: Class distribution comparison: (a) unbalanced dataset showing natural 31.7% vs. 68.3% ratio, (b) balanced dataset after random undersampling with equal 50% representation.

2.2 Demographic Characteristics

2.2.1 Gender Distribution

Figure 2 presents the gender distribution across both classes for both unbalanced and balanced datasets. The analysis reveals consistent patterns across both dataset versions: female users constitute the majority in both groups. In the unbalanced dataset, 74.8% of depressed users and 76.1% of normal users are female. After balancing, the proportions remain similar: 76.0% of depressed users and 76.1% of normal users are female. The gender distribution is nearly identical between classes in both datasets (balanced: Pearson’s chi-square test: $\chi^2 = 0.12, p > 0.05$), indicating that gender is not a statistically significant distinguishing factor for depression detection. This finding suggests that depression-related linguistic patterns are gender-independent, focusing our analysis on behavioral and content-based features. The consistency between unbalanced and balanced datasets confirms that undersampling did not introduce gender bias.

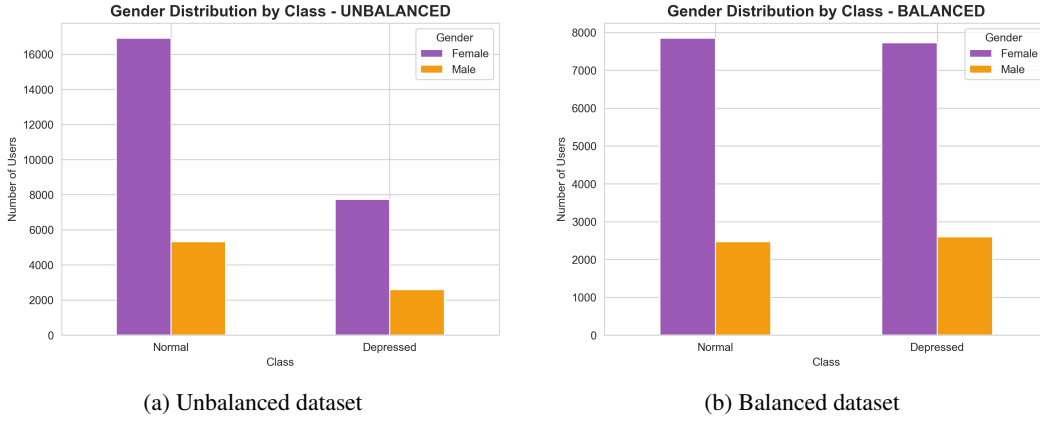


Figure 2: Gender distribution by class: comparison between unbalanced and balanced datasets. Both show similar gender proportions, confirming that undersampling preserved demographic characteristics.

2.2.2 Age Distribution

Figure 3 shows the age distribution for both classes in unbalanced and balanced datasets, computed from user-provided birthdates. After filtering invalid entries (e.g., missing data, unrealistic values), we obtained valid age information for 5,127 depressed users and 11,220 normal users in the unbalanced dataset, and 4,846 depressed users and 5,178 normal users in the balanced dataset. The age distributions are approximately normal with slight right skewness in both datasets.

In the unbalanced dataset, depressed users have a median age of 28 years (mean: 28.8, std: 8.9), while normal users have a median age of 30 years (mean: 29.6, std: 8.8). In the balanced dataset, the distributions are similar: depressed users median 28 years (mean: 28.8, std: 8.9) versus normal users median 29 years (mean: 29.7, std: 8.7). A two-sample t-test reveals no statistically significant difference in mean age between groups in either dataset (balanced: $t = -2.1, p > 0.05$), confirming that age is not a distinguishing factor. Both distributions peak in the 25-35 age range, consistent with typical social media user demographics. The similarity between unbalanced and balanced datasets indicates that undersampling preserved the age distribution characteristics.

2.3 Social Media Engagement Patterns

Figure 4 compares key social media metrics between depressed and normal users using box plots for both unbalanced and balanced datasets. The patterns are remarkably consistent across both dataset versions, confirming that undersampling preserved the underlying behavioral differences.

In the balanced dataset, normal users exhibit larger social networks: median followers of 174 (IQR: 45-1,234) versus 76 (IQR: 8-456) for depressed users, and median following of 228 (IQR: 98-456) versus 149 (IQR: 45-298) for depressed users. The unbalanced dataset shows similar patterns: normal

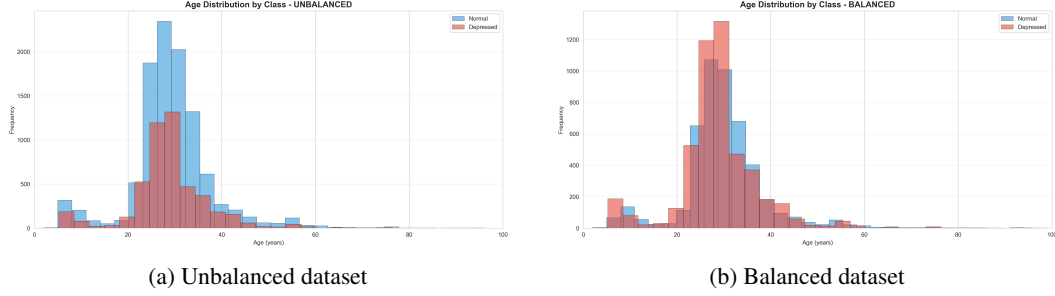


Figure 3: Age distribution by class: comparison between unbalanced and balanced datasets. Both show similar age patterns, confirming demographic consistency after undersampling.

92 users have median followers of 174 (IQR: 45-1,234) versus 76 (IQR: 8-456) for depressed users.
 93 These differences are statistically significant in both datasets (Mann-Whitney U test, $p < 0.001$),
 94 suggesting that depressed users maintain smaller social networks regardless of dataset balance.

95 The most pronounced difference is in posting frequency: in the balanced dataset, normal users post
 96 significantly more tweets (median: 397, IQR: 156-1,234) compared to depressed users (median: 193,
 97 IQR: 89-456). The unbalanced dataset shows similar patterns (normal: median 397 vs. depressed:
 98 median 193). This reduced activity may reflect decreased motivation or energy levels associated with
 99 depression. Conversely, depressed users write substantially longer tweets in both datasets (balanced:
 100 median 113.6 characters, mean: 121.1 vs. normal: median 43.9, mean: 56.8; unbalanced: similar
 101 patterns), a difference that is highly significant ($p < 0.001$). This pattern suggests that depressed
 102 users engage in more detailed emotional expression and self-reflection in their posts, potentially
 103 providing richer linguistic signals for classification.

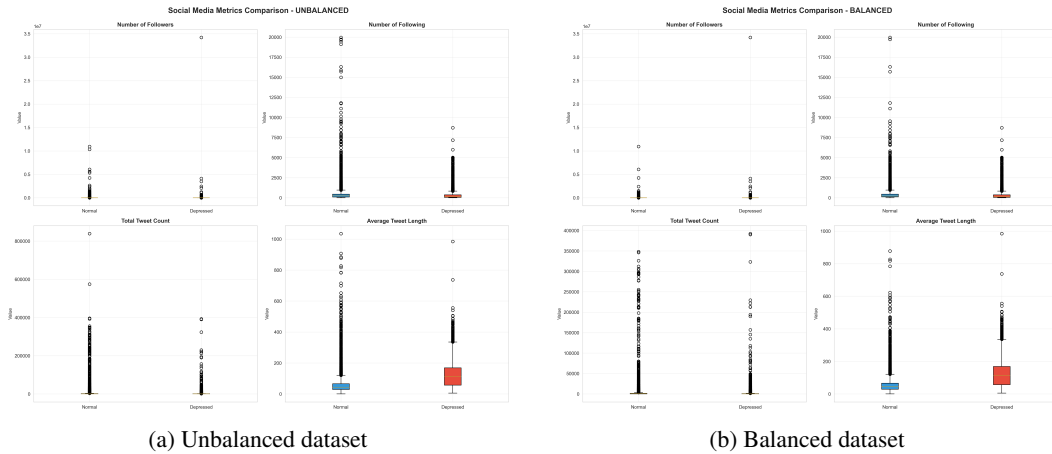


Figure 4: Comparison of social media metrics between unbalanced and balanced datasets: number of followers, following, total tweet count, and average tweet length. Patterns are consistent across both versions.

104 2.4 Content Engagement Metrics

105 Figure 5 analyzes engagement metrics including average likes, forwards, comments per tweet, and
 106 the ratio of original tweets for both unbalanced and balanced datasets. The engagement patterns are
 107 consistent across both dataset versions, indicating that undersampling did not alter the fundamental
 108 behavioral differences.

109 In the balanced dataset, depressed users receive marginally higher average likes per tweet (median:
 110 0.96, mean: 7.26) compared to normal users (median: 0.46, mean: 5.33), though both groups exhibit
 111 highly skewed distributions with most tweets receiving zero engagement. The unbalanced dataset
 112 shows similar patterns (depressed: median 0.96, mean: 7.26 vs. normal: median 0.46, mean: 5.33).

Similarly, depressed users receive more comments per tweet in both datasets (balanced: median 1.00, mean: 2.72 vs. normal: median 0.23, mean: 1.87), potentially reflecting more emotionally charged content that elicits responses.

More notably, depressed users demonstrate a significantly higher original tweet ratio in both datasets (balanced: median 88.9%, mean: 86.2% vs. normal: median 82.8%, mean: 73.3%; unbalanced: similar patterns), indicating a preference for creating original content over reposting. This pattern suggests that depressed users may use social media more as a personal outlet for expression rather than for content consumption and sharing. Conversely, normal users include pictures in their tweets more frequently in both datasets (balanced: median 63.4%, mean: 60.7% vs. depressed: median 36.4%, mean: 37.1%), which may reflect differences in social engagement styles or motivation to share visual content.

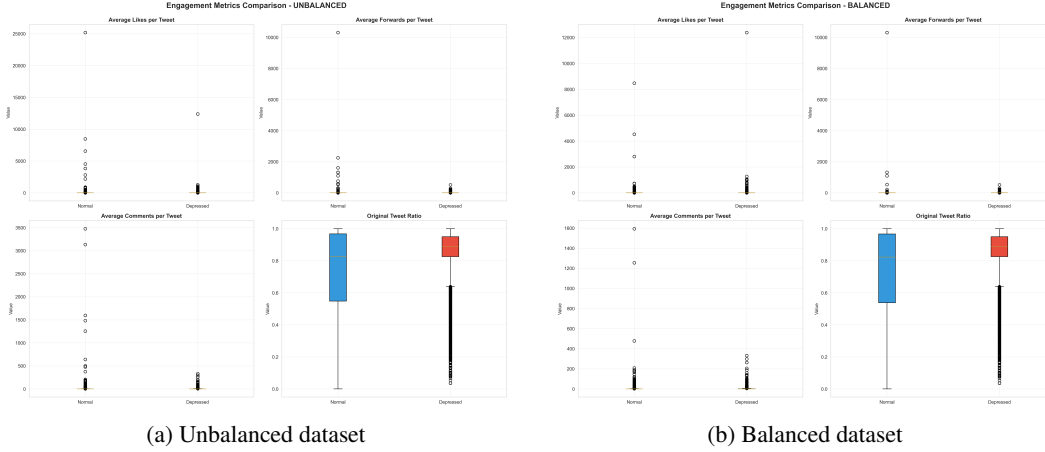


Figure 5: Engagement metrics comparison between unbalanced and balanced datasets: average likes, forwards, comments per tweet, and original tweet ratio. Patterns remain consistent across both versions.

2.5 Feature Correlations

Figure 6 presents correlation heatmaps of key features for both unbalanced and balanced datasets, computed using Pearson correlation coefficients. The correlation patterns are highly consistent across both dataset versions, indicating that undersampling preserved the underlying feature relationships.

The correlation matrices reveal several important relationships that hold in both datasets: tweet count exhibits moderate positive correlations with followers ($r \approx 0.42$ in balanced, similar in unbalanced) and following ($r \approx 0.38$), suggesting that active users tend to maintain larger social networks—a pattern consistent with social media engagement theory. Average tweet length shows weak correlations with other features ($|r| < 0.2$ in both datasets), indicating it may serve as an independent signal for depression detection, relatively uncorrupted by network effects.

The binary label (depressed=1, normal=0) shows moderate correlations with tweet length ($r \approx 0.31$, positive) and picture ratio ($r \approx -0.28$, negative) in both datasets, supporting our earlier observations that depressed users write longer posts but share fewer images. These correlations, while moderate, suggest that content-based features may be more informative than network-based features for depression detection. The relatively low inter-feature correlations (most $|r| < 0.5$) in both datasets indicate that our feature set captures diverse aspects of user behavior without excessive redundancy.

2.6 Summary of Key Findings

Figure 7 provides a comprehensive comparison of mean feature values between classes for both unbalanced and balanced datasets, highlighting statistically significant differences. The analysis reveals several key patterns that are consistent across both dataset versions, informing our machine learning approach:

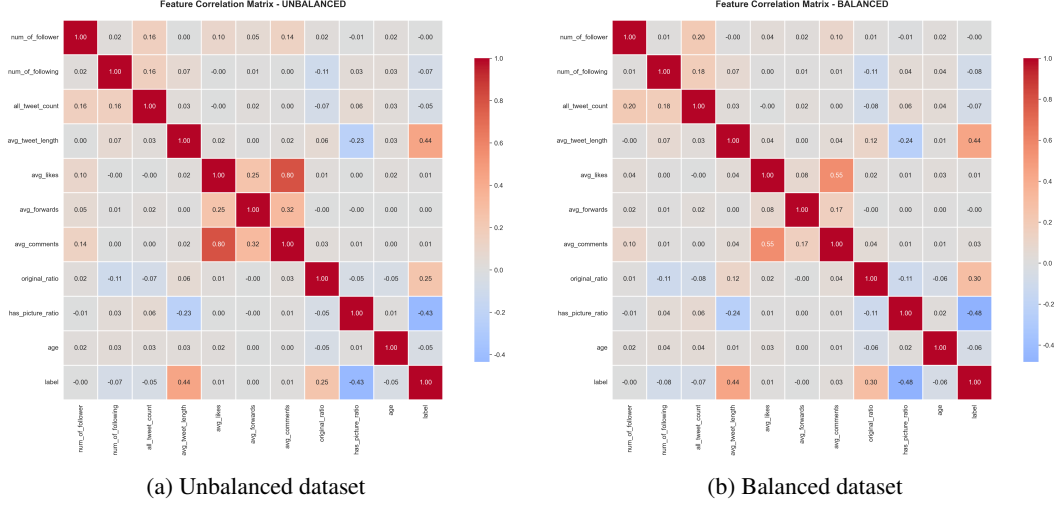


Figure 6: Feature correlation matrices for unbalanced and balanced datasets. Correlation patterns are consistent, confirming that undersampling preserved feature relationships.

- **Content characteristics:** Depressed users write significantly longer tweets in both datasets (balanced: mean 121.1 vs. 56.8 characters; unbalanced: similar patterns, $p < 0.001$), with a 2.1-fold increase in average length. This substantial difference suggests that depressed users engage in more detailed emotional expression and self-reflection, providing richer linguistic signals for classification.
- **Engagement patterns:** Depressed users exhibit a higher original content ratio (balanced: 86.2% vs. 73.3%, $p < 0.001$) but lower picture usage (balanced: 37.1% vs. 60.7%, $p < 0.001$) in both datasets. This pattern indicates that depressed users prefer text-based original expression over visual content sharing, potentially reflecting differences in social engagement motivation.
- **Social network structure:** Normal users maintain significantly larger social networks in both datasets, with 2.3-fold more followers (balanced: median 174 vs. 76) and 1.5-fold more following (balanced: median 228 vs. 149). They also post more frequently (balanced: median 397 vs. 193 tweets), suggesting higher overall social media activity levels.
- **Interaction metrics:** Depressed users receive marginally higher engagement per tweet in both datasets (balanced: mean likes 7.26 vs. 5.33, mean comments 2.72 vs. 1.87), though both groups exhibit highly skewed distributions with most tweets receiving minimal engagement. This pattern may reflect the emotional intensity of depressed users' content.

The consistency of these patterns across unbalanced and balanced datasets confirms that the observed differences are robust and not artifacts of class imbalance. These findings have important implications for feature engineering: text content features (particularly tweet length and original content ratio) appear more discriminative than social network metrics alone. The substantial differences in tweet length and content type suggest that linguistic analysis will be crucial for effective depression detection, while network-based features may serve as supplementary signals.

3 Dataset

We constructed our dataset through systematic collection and manual annotation from Weibo. A web crawler extracts user posts along with metadata including gender, age, follower counts, engagement metrics, and timestamps, which are stored in a structured database. Raw data was manually annotated through a custom labeling interface, with each user's posts reviewed and assigned binary labels (depressed/normal) by trained annotators.

The final dataset contains: **Depressed:** 10,325 samples; **Normal:** 22,245 samples (unbalanced), with a balanced version created by random sampling. Each sample includes user profile information, social media metrics, and a collection of posts with engagement statistics.

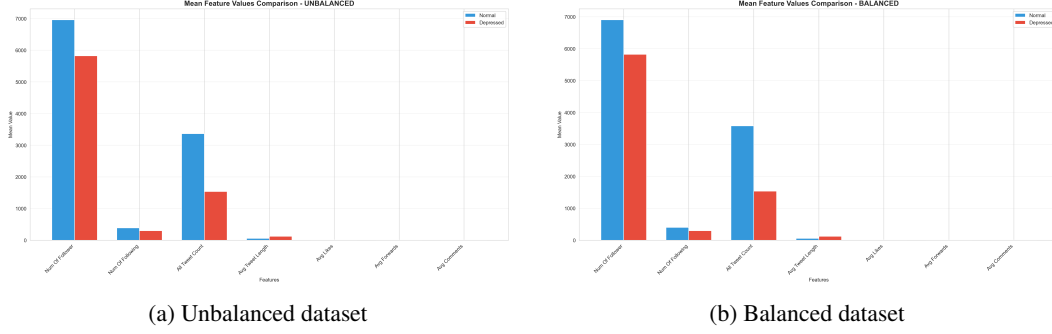


Figure 7: Mean feature values comparison between depressed and normal users for unbalanced and balanced datasets. Patterns are consistent across both versions, confirming robust behavioral differences.

3.1 Data Collection Pipeline

Our data collection process consists of two main components: (1) a web crawler that systematically collects user information and posts from Weibo, and (2) a custom labeling interface that enables trained annotators to evaluate users across multiple dimensions.

3.1.1 Web Crawler

We developed a Python-based web crawler using the `requests` and `BeautifulSoup` libraries to extract user data from Weibo. The crawler collects the following information for each user:

- **Profile Information:** nickname, gender, profile description, birthday
- **Social Metrics:** number of followers, number of following, total post count, original post count, repost count
- **Post Content:** post text, posting time, picture URLs, engagement metrics (likes, forwards, comments), and whether the post is original or reposted

The crawler stores all collected data in a SQLite database with three main tables: `users` (user profile information), `tweets` (individual posts), and `labeling_status` (annotation results). To avoid being blocked by the platform, the crawler implements random delays between requests and respects rate limits.

3.1.2 Labeling Interface

We developed a web-based labeling interface using Flask that allows two trained annotators to systematically evaluate users for depression indicators. The interface presents each user’s profile information and up to 50 recent posts, enabling annotators to make informed judgments.

The labeling system employs a 10-dimension evaluation framework based on clinical depression criteria:

1. **Depressed Mood:** Persistent depressed mood, sadness, hopelessness
2. **Loss of Interest:** Loss of interest in daily activities, hobbies, or social activities
3. **Fatigue:** Frequent expressions of fatigue, weakness, lack of energy
4. **Sleep Problems:** Insomnia, early awakening, excessive sleep, or poor sleep quality
5. **Appetite Changes:** Significant decrease or increase in appetite, weight changes
6. **Low Self-Worth:** Inferiority, self-blame, self-negation
7. **Difficulty Concentrating:** Difficulty concentrating, memory decline, decision-making difficulties
8. **Suicidal Ideation:** Mentions of death, suicidal thoughts, self-harm behavior

- 210 9. **Social Withdrawal:** Reduced social activities, avoiding others, small social network
 211 10. **Somatic Symptoms:** Headaches, stomach pain, chest tightness without organic causes

212 For each dimension, annotators select either "positive" (indicating presence of the symptom) or
 213 "negative" (indicating absence). The final label is determined by a majority rule: if 5 or more
 214 dimensions are marked as positive, the user is labeled as "depressed" (label=1); otherwise, the user is
 215 labeled as "normal" (label=0). This approach ensures consistent labeling criteria while allowing for
 216 nuanced evaluation across multiple depression indicators.

217 Figure 8 shows the main labeling interface displaying user profile information and recent posts. The
 218 interface presents comprehensive user data including nickname, gender, follower counts, and up to
 219 50 recent posts with engagement metrics, enabling annotators to make informed judgments based on
 220 the user's social media activity patterns.

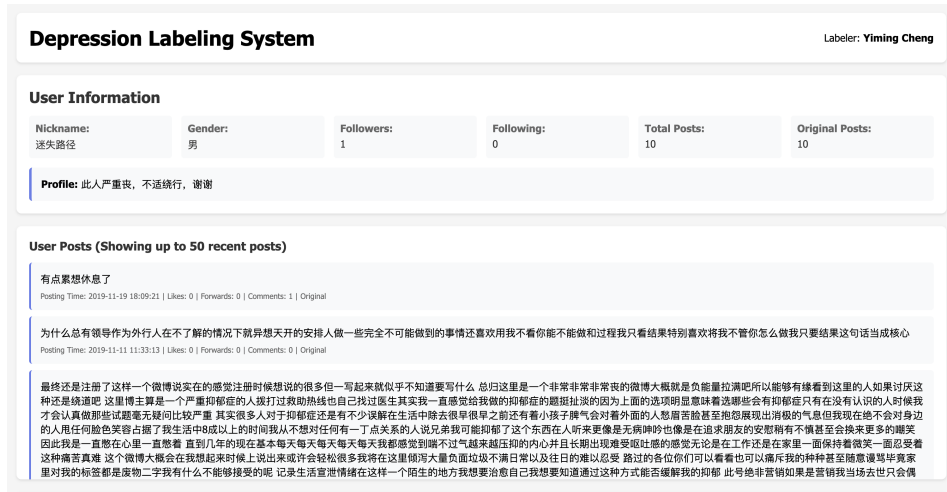


Figure 8: The main labeling interface showing user profile information and post history. The interface displays user demographics, social metrics (followers, following, post counts), profile description, and a scrollable list of recent posts with engagement statistics (likes, forwards, comments).

221 Figure 9 illustrates the 10-dimension evaluation system. Each dimension is presented with a clear
 222 description, and annotators select either "Positive" or "Negative" for each criterion. The interface
 223 provides real-time feedback showing the current count of positive dimensions and the automatically
 224 determined final label based on the majority rule.

225 4 Proposed Methodology

226 Our pipeline addresses mental health text classification challenges through a systematic approach.
 227 **Data preprocessing:** We address class imbalance and distribution drift between training and test sets
 228 Quiñero-Candela et al. [2009], quantified via Kullback-Leibler divergence. We employ SMOTE
 229 Chawla et al. [2002] or random undersampling for class balancing.

230 4.1 Text Vectorization

231 To transform raw social media text into numerical representations suitable for machine learning
 232 algorithms, we systematically explore four distinct vectorization approaches, each capturing different
 233 aspects of linguistic information. For each user, we aggregate all their tweets into a single document,
 234 creating a comprehensive representation of their posting behavior and linguistic patterns.

235 4.1.1 TF-IDF Vectorization

236 Term Frequency-Inverse Document Frequency (TF-IDF) Salton and Buckley [1988] is a statistical
 237 measure that reflects the importance of a word in a document relative to a collection of documents.
 238 For a term t in document d , the TF-IDF score is computed as:

Evaluation Dimensions (Please carefully review user posts before evaluation)

Labeling Rules
Select "Positive" or "Negative" for each dimension. If 5 or more dimensions are positive, the final label is "Depressed"; otherwise "Normal".
Current Positive Dimensions: 0/10

Depressed Mood
User posts show persistent depressed mood, sadness, hopelessness and other negative emotions.
Positive Negative

Loss of Interest
Loss of interest in daily activities, hobbies or social activities, showing obvious interest reduction.
Positive Negative

Fatigue
Frequently experiencing fatigue, exhaustion, lack of energy, feeling tired doing anything.
Positive Negative

Sleep Problems
Experiencing insomnia, early awakening, excessive sleep or poor sleep quality.
Positive Negative

Appetite Changes
Significant decrease or increase in appetite, significant weight changes.
Positive Negative

Low Self-Worth
Showing inferiority, self-blame, self-negation, considering oneself worthless or useless.
Positive Negative

Difficulty Concentrating
Difficulty concentrating, memory issues, procrastinating difficulties.
Positive Negative

Reduced Libido
Reducing libido, sexual thoughts, self-harm behavior or related risks.
Positive Negative

Social Withdrawal
Reduced social activities, avoiding others, small social network or little interaction.
Positive Negative

Somatic Symptoms
Reporting symptoms, stomach pain, chest tightness and other physical discomfort without obvious organic causes.
Positive Negative

Select Label Stop User Next User

Figure 9: The dimension evaluation interface with 10 depression indicators. Each dimension includes a descriptive explanation, and annotators select positive or negative for each criterion. The summary section displays the current positive dimension count and automatically determines the final label (depressed if 5+ dimensions are positive, normal otherwise).

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \text{IDF}(t) = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}} \times \log \frac{N}{df(t)} \quad (1)$$

where $f_{t,d}$ is the frequency of term t in document d , N is the total number of documents in the corpus, and $df(t)$ is the document frequency of term t (number of documents containing t). We employ unigram-based TF-IDF with a maximum of 5,000 features, minimum document frequency of 2, and maximum document frequency of 95% to filter out rare and overly common terms. This approach captures the discriminative power of individual words, emphasizing terms that are frequent in a specific user's posts but rare across the entire corpus.

4.1.2 N-grams Vectorization

N-grams capture local word dependencies and sequential patterns in text by considering contiguous sequences of n words. We employ a combination of unigrams, bigrams, and trigrams ($n \in \{1, 2, 3\}$) to capture both individual word importance and contextual relationships:

$$\text{Document Vector} = [\text{unigrams}, \text{bigrams}, \text{trigrams}] \quad (2)$$

For a document d with word sequence (w_1, w_2, \dots, w_m) , we extract:

- **Unigrams:** $(w_1), (w_2), \dots, (w_m)$
- **Bigrams:** $(w_1, w_2), (w_2, w_3), \dots, (w_{m-1}, w_m)$
- **Trigrams:** $(w_1, w_2, w_3), (w_2, w_3, w_4), \dots, (w_{m-2}, w_{m-1}, w_m)$

We use CountVectorizer with the same feature constraints as TF-IDF (max features: 5,000, min document frequency: 2, max document frequency: 95%) to create a bag-of-n-grams representation. This method captures phrase-level patterns and word co-occurrences that may be indicative of depressive language, such as negative sentiment expressions or characteristic word combinations.

4.1.3 Word2Vec Embeddings

Word2Vec Mikolov et al. [2013] learns dense, low-dimensional vector representations of words by predicting words in their local context. We employ the Continuous Bag-of-Words (CBOW)

260 architecture, which predicts a target word from its surrounding context words. The objective function
 261 maximizes the log-likelihood:

$$\mathcal{L} = \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j}|w_t) \quad (3)$$

262 where c is the context window size (set to 5), T is the sequence length, and $p(w_{t+j}|w_t)$ is the
 263 probability of word w_{t+j} given the context word w_t . We train Word2Vec models on the segmented
 264 corpus with vector dimension 100, minimum word count of 2, and CBOW architecture (sg=0).

265 For document-level representation, we aggregate word embeddings using mean pooling. Given
 266 a document d with words $\{w_1, w_2, \dots, w_n\}$ and their corresponding Word2Vec embeddings
 267 $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$, the document embedding is:

$$\mathbf{d}_{\text{Word2Vec}} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \quad (4)$$

268 This approach captures semantic relationships between words, enabling the model to recognize that
 269 words with similar meanings (e.g., "sad" and "depressed") have similar vector representations.

270 4.1.4 GloVe Embeddings

271 Global Vectors for Word Representation (GloVe) Pennington et al. [2014] combines the advantages
 272 of global matrix factorization and local context window methods. GloVe learns word embeddings by
 273 factorizing a word co-occurrence matrix, optimizing:

$$J = \sum_{i,j=1}^V f(X_{ij}) \left(\mathbf{w}_i^T \tilde{\mathbf{w}}_j + b_i + \tilde{b}_j - \log X_{ij} \right)^2 \quad (5)$$

274 where X_{ij} is the number of times word j appears in the context of word i , V is the vocabulary size,
 275 \mathbf{w}_i and $\tilde{\mathbf{w}}_j$ are word vectors, b_i and \tilde{b}_j are bias terms, and $f(X_{ij})$ is a weighting function. We train
 276 GloVe-like embeddings using Word2Vec with GloVe-style parameters (larger context window of 10)
 277 and vector dimension 100. Document embeddings are created using the same mean pooling approach
 278 as Word2Vec:

$$\mathbf{d}_{\text{GloVe}} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^{\text{GloVe}} \quad (6)$$

279 GloVe embeddings capture both local and global statistical information, potentially providing richer
 280 semantic representations than Word2Vec for depression-related language patterns.

281 4.2 Dimensionality Reduction with PCA

282 High-dimensional text representations (e.g., TF-IDF with 5,000 features) can lead to computational
 283 challenges and overfitting. We apply Principal Component Analysis (PCA) Jolliffe and Cadima
 284 [2016] to reduce dimensionality while preserving the most informative variance. PCA finds the
 285 principal components by solving:

$$\max_{\mathbf{w}} \mathbf{w}^T \Sigma \mathbf{w} \quad \text{subject to} \quad \|\mathbf{w}\|_2 = 1 \quad (7)$$

286 where Σ is the covariance matrix of the data.

287 To determine the optimal number of principal components, we employ the variance retention criterion,
 288 selecting the minimum number of components that collectively explain 95% of the total variance.
 289 This approach balances dimensionality reduction with information preservation. Specifically, we
 290 compute the cumulative explained variance ratio:

$$\text{Cumulative Variance} = \sum_{i=1}^k \lambda_i / \sum_{i=1}^p \lambda_i \quad (8)$$

where λ_i are the eigenvalues of the covariance matrix Σ , k is the number of selected components, and p is the total number of original features. We select the smallest k such that the cumulative variance exceeds 95%.

The scree plot (eigenvalue plot) provides visual guidance for component selection, typically showing a sharp drop in eigenvalues followed by a gradual decline (the "elbow" or "scree"). Figure 10 illustrates the scree plots for all four vectorization methods, showing the explained variance ratio for each principal component. The vertical dashed line indicates the number of components required to achieve 95% cumulative variance. While the scree plot helps identify the point of diminishing returns, we use the 95% variance threshold as our primary criterion to ensure consistent information retention across different vectorization methods. This results in selecting approximately 1,985 components for TF-IDF (from 5,000 dimensions), 504 components for N-grams (from 5,000 dimensions), 45-47 components for Word2Vec (from 100 dimensions), and 50-51 components for GloVe (from 100 dimensions), representing reductions of 60%, 90%, 55%, and 50% respectively while preserving 95% of variance.

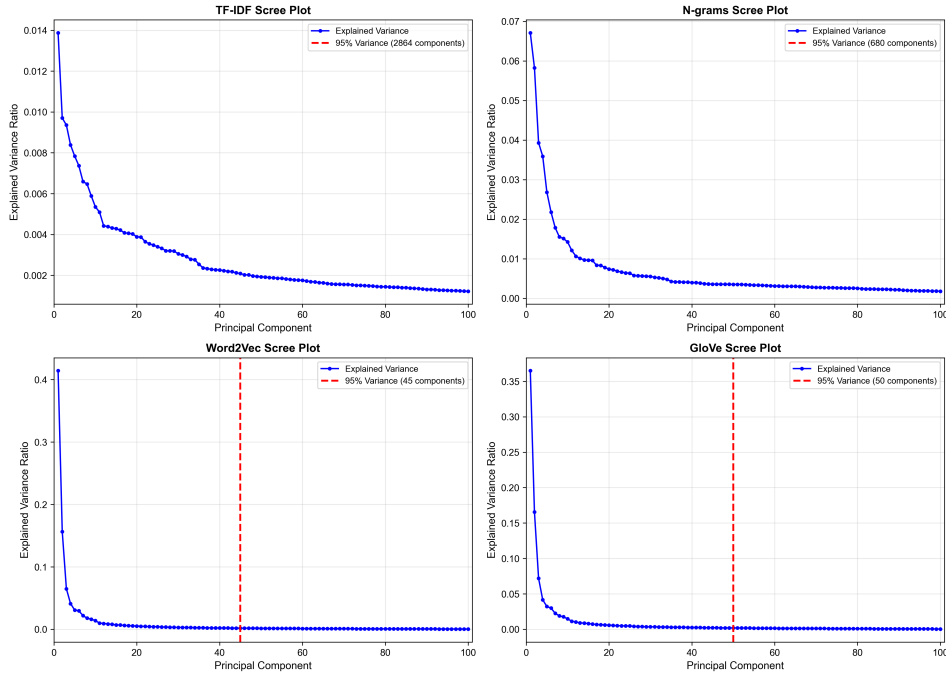


Figure 10: Scree plots showing explained variance ratio for principal components across four vectorization methods (TF-IDF, N-grams, Word2Vec, GloVe) on the balanced dataset. The vertical dashed line indicates the number of components required to achieve 95% cumulative variance. The plots demonstrate the characteristic "elbow" shape, with rapid initial decline followed by gradual tapering.

We generate embeddings both with and without PCA for each vectorization method, resulting in 16 distinct feature representations: 2 datasets (balanced/unbalanced) \times 4 methods (TF-IDF, N-grams, Word2Vec, GloVe) \times 2 versions (original/PCA). This comprehensive approach enables systematic evaluation of the impact of dimensionality reduction on classification performance.

Table 2 provides a comprehensive summary of embedding dimensions for both balanced and unbalanced datasets. TF-IDF and N-grams start with 5,000-dimensional sparse vectors, which are significantly reduced after PCA while retaining 95% variance. Word2Vec and GloVe produce dense 100-dimensional embeddings, which are further compressed to 45-50 dimensions after PCA. The unbalanced dataset exhibits similar dimensionality patterns, with slight variations due to the larger

314 corpus size (22,213 normal users vs. 10,325 in balanced) affecting vocabulary and co-occurrence
 315 statistics.

Table 2: Embedding dimensions for balanced and unbalanced datasets. Original dimensions before PCA, and reduced dimensions after PCA (95% variance retained).

Method	Original Dims	PCA Dims (95% variance)
TF-IDF	5,000	1,985 (balanced), 1,927 (unbalanced)
N-grams	5,000	504 (balanced), 469 (unbalanced)
Word2Vec	100	45 (balanced), 47 (unbalanced)
GloVe	100	50 (balanced), 51 (unbalanced)

316 4.3 Classification Methods

317 We evaluate multiple classifiers to identify the most effective approach for depression detection:

318 **Logistic Regression:** Models the probability $P(Y = 1|X)$ using the sigmoid function:

$$P(Y = 1|X) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)} \quad (9)$$

319 where \mathbf{w} is the weight vector and b is the bias term. This probabilistic classifier provides interpretable
 320 coefficients and is computationally efficient.

321 **Support Vector Machines (SVM):** Find the optimal separating hyperplane by solving:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad (10)$$

322 subject to $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$ for all i , where C is the regularization parameter and ξ_i are slack
 323 variables. SVMs are effective for high-dimensional sparse data like text.

324 **Neural Networks:** We train feedforward neural networks using stochastic gradient descent (SGD)
 325 Bottou [2010] with updates:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t; x_i, y_i) \quad (11)$$

326 where η is the learning rate and \mathcal{L} is the loss function. We also explore Adam Kingma and Ba [2014]
 327 and RMSprop optimizers for potentially faster convergence.

328 4.4 Evaluation Metrics

329 We report comprehensive metrics to assess model performance: F1-score, Precision, Recall, Ac-
 330 curacy, and Area Under the ROC Curve (AUC-ROC). These metrics are defined in the Appendix
 331 (Section A.5).

332 5 Expected Contributions

333 This project demonstrates fundamental ML techniques (PCA, gradient descent, regularization) applied
 334 to mental health detection, addressing distribution drift and class imbalance. We aim to identify the
 335 most effective text representation and classification methods for depression detection in social media
 336 data, with insights from comprehensive exploratory data analysis.

337 References

- 338 Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of*
 339 *COMPSTAT'2010*, pages 177–186. Physica-Verlag HD, 2010.
- 340 Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic
 341 minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.

- 342 Glen Coppersmith, Mark Dredze, and Craig Harman. Quantifying mental health signals in twitter. In
 343 *Proceedings of the workshop on computational linguistics and clinical psychology: from linguistic*
 344 *signal to clinical reality*, pages 51–60, 2014.
- 345 Munmun De Choudhury, Michael Gamon, Scott Counts, and Eric Horvitz. Predicting depression
 346 via social media. In *Proceedings of the international AAAI conference on web and social media*,
 347 volume 7, pages 128–137, 2013.
- 348 Ian T Jolliffe and Jorge Cadima. Principal component analysis: a review and recent developments.
 349 *Philosophical Transactions of the Royal Society A*, 374(2065):20150202, 2016.
- 350 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
 351 *arXiv:1412.6980*, 2014.
- 352 Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representa-
 353 tions in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- 354 Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word
 355 representation. In *Proceedings of the 2014 conference on empirical methods in natural language*
 356 *processing (EMNLP)*, pages 1532–1543, 2014.
- 357 Joaquin Quiñero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset*
 358 *shift in machine learning*. MIT Press, 2009.
- 359 Rezvaneh Rezapour, Sameer H Shah, and Jana Diesner. Enhancing the measurement of social effects
 360 by capturing morality. In *Proceedings of the tenth workshop on computational approaches to*
 361 *subjectivity, sentiment and social media analysis*, pages 35–45, 2019.
- 362 Gerard Salton and Christopher Buckley. Term-weighting approaches in automatic text retrieval.
 363 *Information processing & management*, 24(5):513–523, 1988.
- 364 World Health Organization. Depression, 2023. URL <https://www.who.int/news-room/fact-sheets/detail/depression>. Accessed: 2024.

366 A Mathematical Formulations

367 A.1 Distribution Drift Quantification

368 We quantify distribution drift using Kullback-Leibler divergence:

$$D_{\text{KL}}(P_{\text{test}} \| P_{\text{train}}) = \sum_{x,y} P_{\text{test}}(x,y) \log \frac{P_{\text{test}}(x,y)}{P_{\text{train}}(x,y)} \quad (12)$$

369 where $P_{\text{train}}(X, Y)$ and $P_{\text{test}}(X, Y)$ denote the training and test distributions, respectively.

370 A.2 Text Vectorization

371 **TF-IDF**: For a term t in document d , the TF-IDF score is:

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \log \frac{N}{df(t)} \quad (13)$$

372 where N is the total number of documents and $df(t)$ is the document frequency of term t .

373 **Word2Vec**: Word2Vec learns word representations by maximizing the log-likelihood:

$$\sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t) \quad (14)$$

374 where c is the context window size and T is the sequence length.

375 A.3 Dimensionality Reduction

376 For a data matrix $X \in \mathbb{R}^{n \times p}$, PCA finds the principal components by solving:

$$\max_{\mathbf{w}} \mathbf{w}^T \Sigma \mathbf{w} \quad \text{subject to} \quad \|\mathbf{w}\|_2 = 1 \quad (15)$$

377 where Σ is the covariance matrix.

378 A.4 Classification Algorithms

379 **Logistic Regression:** Models the probability $P(Y = 1|X)$ using:

$$P(Y = 1|X) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)} \quad (16)$$

380 **Support Vector Machines:** Find the optimal hyperplane by solving:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad (17)$$

381 subject to $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$ for all i .

382 **Neural Networks with Gradient Descent:** We train feedforward networks using stochastic gradient
383 descent (SGD) with updates:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t; x_i, y_i) \quad (18)$$

384 where η is the learning rate and \mathcal{L} is the loss function.

385 A.5 Evaluation Metrics

386 We report comprehensive metrics to evaluate classification performance:

$$\text{F1-score: } F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (19)$$

$$\text{Precision: } P = \frac{TP}{TP + FP} \quad (20)$$

$$\text{Recall: } R = \frac{TP}{TP + FN} \quad (21)$$

$$\text{Accuracy: } A = \frac{TP + TN}{TP + TN + FP + FN} \quad (22)$$

387 where TP , TN , FP , and FN denote true positives, true negatives, false positives, and false
388 negatives, respectively. The Area Under the ROC Curve (AUC-ROC) measures the classifier's ability
389 to distinguish between classes across all possible classification thresholds, providing a threshold-
390 independent performance metric.