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

Yiming Cheng¹ Yi Wu¹

¹Department of Computer Science, University of Chicago, MPCS Predoc
eaminchan@uchicago.edu, yiwu@uchicago.edu

Abstract

We propose a machine learning framework for depression detection from social media text, addressing distribution drift. We employ multiple vectorization techniques (TF-IDF, N-grams, Word2Vec, GloVe), PCA, and evaluate classifiers (logistic regression, SVM, neural networks with SGD). Methodology includes preprocessing, SMOTE for class balancing, and evaluation using F1-score, precision, recall, AUC-ROC. Dataset: 10,325 depressed and 22,245 normal samples from X (Twitter).

1 Introduction and problem statement

Depression affects millions worldwide; early detection is crucial World Health Organization [2023]. Social media provides rich linguistic data revealing psychological states De Choudhury et al. [2013]. Previous work demonstrates feasibility of depression detection from social media Coppersmith et al. [2014], Rezapour et al. [2019]. We propose an ML system analyzing social media posts to identify at-risk individuals, addressing distribution drift challenges.

2 Dataset

We constructed a dataset via systematic collection and manual annotation from X (Twitter) through web scraping. A crawler extracts user posts with metadata (gender, age, follower counts, engagement metrics, timestamps) stored in a database. Raw data was manually annotated through a custom labeling interface, with each user's posts reviewed and assigned binary labels (depressed/normal). Final dataset: **Depressed**: 10,325 samples; **Normal**: 22,245 samples.

3 Proposed methodology

Our pipeline addresses mental health text classification challenges. **Data preprocessing**: We address class imbalance and distribution drift between training and test sets Quiñonero-Candela et al. [2009], quantified via Kullback-Leibler divergence (see Appendix). We employ SMOTE Chawla et al. [2002] or random undersampling for class balancing.

Text vectorization: We explore multiple representations: (1) **TF-IDF** Salton and Buckley [1988] for term weighting; (2) **N-grams** (unigrams, bigrams, trigrams) to capture local word dependencies; (3) **Word embeddings** using pre-trained Word2Vec Mikolov et al. [2013] and GloVe Pennington et al. [2014] for semantic representations. Detailed formulations are provided in the Appendix.

Dimensionality reduction: We apply Principal Component Analysis (PCA) Jolliffe and Cadima [2016] to reduce high-dimensional text features, retaining components explaining 95% variance. The optimization formulation is detailed in the Appendix.

Classification: We evaluate multiple classifiers: (1) **Logistic Regression** for probabilistic classification; (2) **Support Vector Machines (SVM)** for maximum-margin classification; (3) **Neural Networks** trained via stochastic gradient descent (SGD) Bottou [2010], also exploring Adam Kingma and Ba [2014] and RMSprop optimizers. Mathematical formulations are provided in the Appendix.

Evaluation: We report comprehensive metrics including F1-score, Precision, Recall, Accuracy, and Area Under the ROC Curve (AUC-ROC). Detailed metric definitions are provided in the Appendix.

4 Expected contributions

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

41 **A Mathematical Formulations**

42 **A.1 Distribution Drift Quantification**

43 We quantify distribution drift using Kullback-Leibler divergence:

$$D_{\text{KL}}(P_{\text{test}} \| P_{\text{train}}) = \sum_{x,y} P_{\text{test}}(x,y) \log \frac{P_{\text{test}}(x,y)}{P_{\text{train}}(x,y)} \quad (1)$$

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

45 **A.2 Text Vectorization**

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

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \log \frac{N}{df(t)} \quad (2)$$

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

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

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

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

50 **A.3 Dimensionality Reduction**

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

$$\max_{\mathbf{w}} \mathbf{w}^T \Sigma \mathbf{w} \quad \text{subject to} \quad \|\mathbf{w}\|_2 = 1 \quad (4)$$

52 where Σ is the covariance matrix.

53 **A.4 Classification Algorithms**

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

$$P(Y = 1|X) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)} \quad (5)$$

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

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad (6)$$

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

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

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t; \mathbf{x}_i, y_i) \quad (7)$$

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

60 **A.5 Evaluation Metrics**

61 We report comprehensive metrics:

$$\text{F1-score: } F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (8)$$

$$\text{Precision: } P = \frac{TP}{TP + FP} \quad (9)$$

$$\text{Recall: } R = \frac{TP}{TP + FN} \quad (10)$$

$$\text{Accuracy: } A = \frac{TP + TN}{TP + TN + FP + FN} \quad (11)$$

62 where TP , TN , FP , and FN denote true positives, true negatives, false positives, and false negatives, respectively.

64 **References**

- 65 Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of
66 COMPSTAT'2010*, pages 177–186. Physica-Verlag HD, 2010.
- 67 Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic
68 minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- 69 Glen Coppersmith, Mark Dredze, and Craig Harman. Quantifying mental health signals in twitter. In
70 *Proceedings of the workshop on computational linguistics and clinical psychology: from linguistic
71 signal to clinical reality*, pages 51–60, 2014.
- 72 Munmun De Choudhury, Michael Gamon, Scott Counts, and Eric Horvitz. Predicting depression
73 via social media. In *Proceedings of the international AAAI conference on web and social media*,
74 volume 7, pages 128–137, 2013.
- 75 Ian T Jolliffe and Jorge Cadima. Principal component analysis: a review and recent developments.
76 *Philosophical Transactions of the Royal Society A*, 374(2065):20150202, 2016.
- 77 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
78 arXiv:1412.6980*, 2014.
- 79 Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representa-
80 tions in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- 81 Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word
82 representation. In *Proceedings of the 2014 conference on empirical methods in natural language
83 processing (EMNLP)*, pages 1532–1543, 2014.
- 84 Joaquin Quiñonero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset
85 shift in machine learning*. MIT Press, 2009.
- 86 Rezvaneh Rezapour, Sameer H Shah, and Jana Diesner. Enhancing the measurement of social effects
87 by capturing morality. In *Proceedings of the tenth workshop on computational approaches to
88 subjectivity, sentiment and social media analysis*, pages 35–45, 2019.
- 89 Gerard Salton and Christopher Buckley. Term-weighting approaches in automatic text retrieval.
90 *Information processing & management*, 24(5):513–523, 1988.
- 91 World Health Organization. Depression, 2023. URL [https://www.who.int/news-room/
92 fact-sheets/detail/depression](https://www.who.int/news-room/fact-sheets/detail/depression). Accessed: 2024.