
Predicting Depression from Synthetic Mental Health Survey Data with Classical and Neural Models

COMP 560 Final Project

Tejaswi Paladugu, Nidhi Padala, Harjas Kaur, Sanjana Nalla

Problem

- Depression
 - Prevalent mental health condition, particularly impacting young adults and students, affecting daily life, academic performance, and long term wellbeing
- Our Goal
 - To use models to systematically examine how demographic, academic, and lifestyle factors are associated with depression
- Kaggle Task
 - Classification: Predict the depression label given the features

Models

- Gaussian Naive Bayes
 - A sanity-check to establish a baseline performance and test the initial data pipeline
- Logistic Regression with One-Hot Encoding
 - A more informative linear model that provides interpretable coefficients (feature importance) for categorical variables
- Unsupervised Learning using PCA (with two components)
 - To explore dimensionality reduction and understand the separation in the feature space
- Neural Network
 - A model to capture potential nonlinear interactions between features, aiming for the highest predictive performance

Models

	Gaussian Naive Bayes	Logistic Regression with One-Hot Encoding	Unsupervised Learning with PCA (2 components)	Neutral Networks
Training Accuracy	0.8589 ± 0.0022	0.9387 ± 0.0010	0.9901	0.9399
Training AUC-ROC	0.9235 ± 0.0018	0.9744 ± 0.0009	0.9993	0.9763
Notes	Simple baseline, lowest performance expected	Coefficients show feature associations, outperforms Naive Bayes	Overfitting on only 2 components, should not be overinterpreted	Most flexible, captures nonlinear relationships

Analysis

- Gaussian Naive Bayes
 - Lowest accuracy (≈ 0.86) and AUC (≈ 0.92)- expected because this is a baseline
 - Indicates depression in this dataset depends on interactions between factors and not isolated effects
- Logistic Regression (One-Hot)
 - Strong linear performer (≈ 0.94 accuracy, ≈ 0.97 AUC)
 - One-hot encoding allows it to capture detailed demographic, academic, and lifestyle effects
 - Features like gender, city, job/student status, and satisfaction variables contribute to predicting depression
- PCA-Based Classifier (2 Components)
 - Extremely high performance (≈ 0.99 accuracy, ≈ 0.999 AUC) from only two components shows the dataset has strong global structure
 - Much of the information needed to predict depression lies in a 2D linear projection of the numeric features.
 - Suggests the dataset is highly separable and dominated by a few major factors
- Neural Network
 - Captures nonlinear patterns but does not outperform linear models
 - Provides probability confidence, showing most depressed cases cluster near high predicted probabilities and non-depressed cases near zero
 - Clear separation in predicted probabilities
- Overall Insights
 - Depression risk in this dataset is highly predictable with even simple or linear models