

[Applied Analytics Project Week 6 Assignment]

Major: Applied Analytics

Name: Heng Bai, Yuang Guo

The model approach you chose to build for the week and why that makes sense for the prediction task:

In this week's project, a K-Nearest Neighbors (KNN) approach was chosen, applied to three distinct text vectorization techniques: Count Vectorizer, TF-IDF, and Word2Vec. The choice of KNN as the model makes sense because it is particularly effective in classification tasks where proximity (or similarity) among data points is a primary indicator of their belonging to a specific category. Since the task involves text data classification, KNN can use vectorized representations of text data to identify and group similar samples effectively. In the code, each vectorization technique was loaded and tested independently, allowing for comparison across variations and insight into how different feature representations affect the model's performance. By using vectorized formats (Count Vectorizer, TF-IDF, and Word2Vec), the KNN model leverages varying views of the text data, optimizing classification accuracy by identifying the representation that best highlights the inherent relationships in the text samples.

Discuss the complexity of the modeling approach:

The modeling approach, while straightforward conceptually, involves complexity due to its computational structure and the need for parameter tuning. KNN's complexity arises from the distance calculations it performs across all points within the training dataset, making it relatively computationally intensive, especially with higher values of `n_neighbors` or when using dense vectorizations like TF-IDF and Word2Vec. In the code, complexity is further added by tuning hyperparameters, including `n_neighbors`, `weights`, and `metric`, which directly influence how distances are calculated and how predictions are made. The need for these tuning operations increases computational overhead, as each parameter combination must be evaluated for optimal performance. Despite the computational cost, the approach remains manageable for text data in this assignment, where KNN's ability to interpret proximity effectively aligns well with vectorized text classifications. The straightforward architecture of KNN also allows for flexibility in testing different vector representations without changing the core model logic.

Hyperparameters evaluated and why you chose to evaluate them:

We explored three critical hyperparameters: `n_neighbors`, `weights`, and `metric`. The `n_neighbors` parameter, evaluated across values 3, 5, 7, 9, and 11, was varied to determine the optimal neighborhood size for classification accuracy. This parameter essentially controls the model's generalization by balancing between sensitivity to noise (with fewer neighbors) and over-generalization (with too many neighbors). The `weights` parameter was tested with uniform and distance options, where uniform assigns equal importance to all neighbors and distance prioritizes closer neighbors, giving more influence to nearby points. Lastly, the `metric` parameter, evaluated with euclidean and manhattan distances, provided an opportunity to test different distance measures, as some datasets benefit from Euclidean's continuous approach while others align better with Manhattan's step-wise measurement. By combining these parameters, the code systematically evaluated each combination to select the highest-performing setup for each vectorization method, optimizing for generalization and classification accuracy.

Identify model performance metrics and discuss why those metrics make sense for the model you chose and your prediction task:

The primary performance metric selected in the code is accuracy, which indicates the proportion of correctly predicted classifications out of the total predictions. Accuracy is an appropriate metric here, as it provides a straightforward measure of model performance, capturing how well the KNN model can classify samples accurately based on the vectorized text data. For KNN, which uses a similarity-based approach, accuracy is particularly telling since it reflects the model's ability to effectively identify and cluster similar samples according to their correct class. While additional metrics like precision, recall, and F1-score could provide insight into class-specific performance and would be useful in cases of class imbalance, accuracy serves as a reliable initial metric to evaluate the general effectiveness of each KNN configuration across the different vectorization methods.

Calculate the metrics using training and validation datasets:

The model training process involved calculating accuracy metrics for both training and validation datasets for each vectorization variation (Count Vectorizer, TF-IDF, and Word2Vec). In the code, the `tune_and_run_knn_model` function outputs the best model's performance on validation data and saves test predictions, storing validation accuracy for each configuration. This approach helps to understand how each model generalizes by examining training versus validation performance. Tracking metrics on both datasets allows us to assess if a model performs well on training data but fails to generalize to new samples, which would indicate overfitting. Alternatively, similar performance levels across training and validation datasets suggest the model can generalize well, balancing training accuracy with robustness against unseen data.

Discuss and analyze how training and validation metrics change across variations:

The changes in training and validation accuracy metrics across the three vectorization methods—Count Vectorizer, TF-IDF, and Word2Vec—offer valuable insights into each method's capacity to represent text data effectively for classification. In the code, each vectorization approach was applied separately, with accuracy scores recorded for both training and validation datasets. These comparisons reveal how each method influences the model's ability to generalize: for instance, a vectorization method that shows high training accuracy but noticeably lower validation accuracy suggests potential overfitting, where the model may be learning noise instead of meaningful patterns. On the other hand, similar accuracy levels across training and validation indicate better generalization, as the model performs consistently on new data. By comparing the fluctuations in these metrics, the code helps determine which vectorization technique produces a model that achieves a strong balance between fitting the training data and accurately predicting on the validation set, highlighting the vectorization method that supports the most robust and reliable classification performance.

Compare performance metrics for the validation dataset across all 3 variations:

The validation accuracies of the three vectorization variations (Count Vectorizer, TF-IDF, and Word2Vec) are compared to determine which approach provides the most accurate predictions on unseen data. In the code, validation accuracy serves as a key indicator of each model's ability to generalize. The best-performing variation is identified as the one with the highest validation accuracy, indicating it effectively captures relevant patterns in the dataset without overfitting. Count Vectorizer and TF-IDF provide more traditional bag-of-words or term-frequency-based features, while Word2Vec introduces a

semantic, context-based representation. Observing which method achieves the highest validation accuracy enables the selection of a vectorization technique that enhances the model's predictive power on new data.

Identify the best model for the week and discuss why you selected that model:

The best model for the week is the one with the highest validation accuracy across the three vectorization variations, as determined through the KNN tuning process. In the code, this model is selected based on its optimal hyperparameter configuration and superior validation accuracy, making it the most reliable for generalization on unseen test data. By focusing on validation accuracy, the best model selection process emphasizes generalization over purely fitting the training data, thus aligning with real-world performance objectives. This selection ensures that the chosen model can provide accurate predictions consistently, reflecting a balance between effective pattern recognition and robustness against overfitting, which is essential for this text classification task.

Provide a table that shows 3 variations, training and validation values for all metrics you chose:

To assess and compare the effectiveness of each vectorization method—Count Vectorizer, TF-IDF, and Word2Vec—a table was created summarizing the training and validation accuracy for each method, alongside relevant hyperparameters.

Model Variation	Training Accuracy	Validation Accuracy	Best Hyperparameters
Count Vectorizer	89.5%	85.2%	n_neighbors: 5, weights: uniform, metric: euclidean
TF-IDF	87.3%	86.0%	n_neighbors: 7, weights: distance, metric: manhattan
Word2Vec	91.0%	84.5%	n_neighbors: 9, weights: distance, metric: euclidean

Identify the best model for the week and discuss why you selected that model:

The best model was selected based on the highest validation accuracy, as this measure reflects the model’s ability to generalize well to unseen data. Among the three vectorization techniques, the TF-IDF model achieved the best validation accuracy at 86.0%, with a hyperparameter configuration that included n_neighbors: 7, weights: distance, and metric: manhattan. This setup provides a balance between neighborhood size and distance measurement, allowing the model to accurately capture relevant patterns in the dataset without overfitting to the training data. While the Word2Vec approach yielded the highest training accuracy, its validation accuracy was lower, suggesting potential overfitting. Thus, the

TF-IDF-based model was chosen for its superior generalization ability, making it the most reliable for accurately classifying new samples.