

# **UE23CS352A: Machine Learning**

## **Lab Week 12: Naive Bayes Classifier**

Name: NISHA SURESH PATIL

SRN: PES2UG23CS393

Date: 30-10-2025

### **Introduction**

The goal of this lab is to learn how to classify biomedical abstract sentences into specific categories like BACKGROUND, METHODS, RESULTS, OBJECTIVE, and CONCLUSION using probabilistic approaches. You'll start by implementing a Multinomial Naive Bayes classifier from scratch, gaining insight into how the model works under the hood. Then, you'll leverage scikit-learn to vectorize text, build models, and improve performance through hyperparameter tuning with GridSearchCV. Finally, you'll explore ensemble learning by approximating a Bayes Optimal Classifier (BOC) using a soft voting classifier that combines multiple diverse models. Altogether, these exercises help build an understanding of probabilistic reasoning, model optimization, and how ensembles can enhance text classification.

### **Methodology**

The implementation started with building a Multinomial Naive Bayes (MNB) classifier from scratch to gain a deeper understanding of the underlying probabilistic principles. Using CountVectorizer features, the model calculated log prior probabilities and log likelihoods for each word-class pair, applying Laplace smoothing to handle words that didn't appear in the training data. During prediction, the model summed the log probabilities across all features to determine the most likely class for each sentence.

Next, the Bayes Optimal Classifier (BOC) was approximated by creating an ensemble of five diverse models: MultinomialNB, Logistic Regression, Random Forest, Decision Tree, and K-Nearest Neighbors. Each model was trained on sampled data, and their posterior probabilities were combined using a soft voting approach, effectively weighting each model according to its performance. This strategy allowed the ensemble to make more robust and accurate predictions, approximating optimal Bayesian decision-making.

Results and Analysis (Screenshots of plots and metrics):

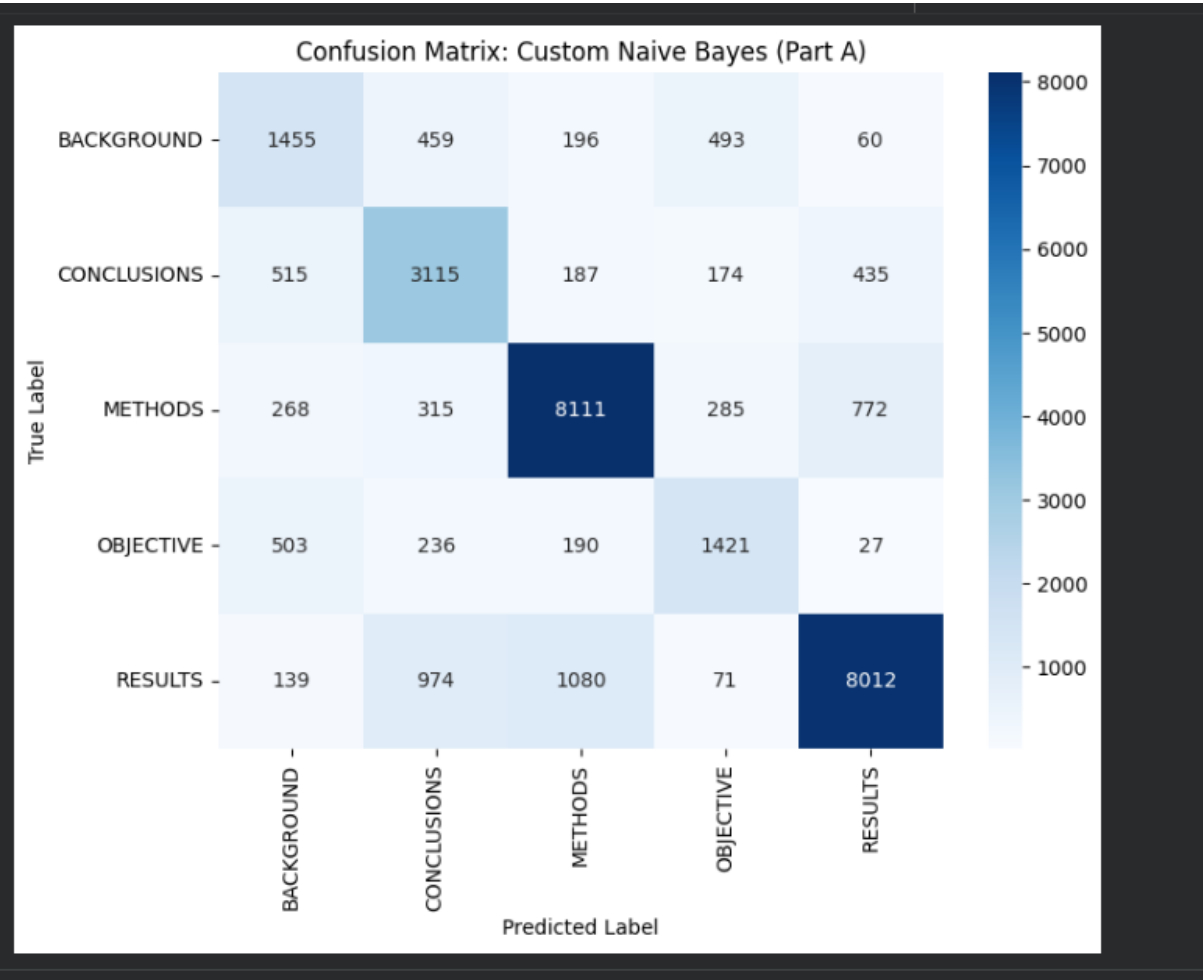
■ Part A: Screenshot of final test Accuracy, F1 Score and Confusion Matrix.

2x

=== Test Set Evaluation (Custom Count-Based Naive Bayes) ===  
Accuracy: 0.7498

	precision	recall	f1-score	support
BACKGROUND	0.51	0.55	0.52	2663
CONCLUSIONS	0.61	0.70	0.65	4426
METHODS	0.83	0.83	0.83	9751
OBJECTIVE	0.58	0.60	0.59	2377
RESULTS	0.86	0.78	0.82	10276
accuracy			0.75	29493
macro avg	0.68	0.69	0.68	29493
weighted avg	0.76	0.75	0.75	29493

Macro-averaged F1 score: 0.6836



## ■ Part B: Screenshot of best hyperparameters found and their resulting F1 score.

```
Training initial Naive Bayes pipeline...
Training complete.

=== Test Set Evaluation (Initial Sklearn Model) ===
Accuracy: 0.7612
      precision    recall  f1-score   support

 BACKGROUND      0.66      0.37      0.47      2663
 CONCLUSIONS  0.64      0.70      0.67      4426
      METHODS      0.79      0.87      0.83      9751
 OBJECTIVE      0.73      0.39      0.51      2377
      RESULTS      0.81      0.87      0.84     10276

 accuracy          0.76          0.76          0.75          29493
 macro avg         0.73          0.64          0.66          29493
 weighted avg      0.76          0.76          0.75          29493

Macro-averaged F1 score: 0.6631

Starting Hyperparameter Tuning on Development Set...
Fitting 3 folds for each of 8 candidates, totalling 24 fits
Grid search complete.

Best Cross-Validation Score (Macro F1): 0.6456
Best Parameters Found: {'nb__alpha': 0.1, 'tfidf__ngram_range': (1, 2)}
```

## ■ Part C:

### 1. Screenshot of SRN and sample size.

```
Please enter your full SRN (e.g., PE51UG22CS345): PE52UG23CS393
Using dynamic sample size: 10393
Actual sampled training set size used: 10393
NaiveBayes trained successfully.
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be removed in a future version.
  warnings.warn(
LogisticRegression trained successfully.
RandomForest trained successfully.
DecisionTree trained successfully.
KNN trained successfully.

Training all base models...
All base models trained.

Calculating posterior weights from validation performance...
NaiveBayes Validation F1: 0.6090
/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:1247: FutureWarning: 'multi_class' was deprecated in version 1.5 and will be removed in a future version.
  warnings.warn(
LogisticRegression Validation F1: 0.5842
RandomForest Validation F1: 0.5156
DecisionTree Validation F1: 0.3533
KNN Validation F1: 0.1828

Posterior Weights (normalized):
NaiveBayes: 0.232
LogisticRegression: 0.226
RandomForest: 0.211
DecisionTree: 0.180
KNN: 0.151

Fitting the VotingClassifier (BOC approximation)...
Fitting complete.

Predicting on test set...
```

### 2. Screenshot of BOC final Accuracy, F1 Score and Confusion Matrix

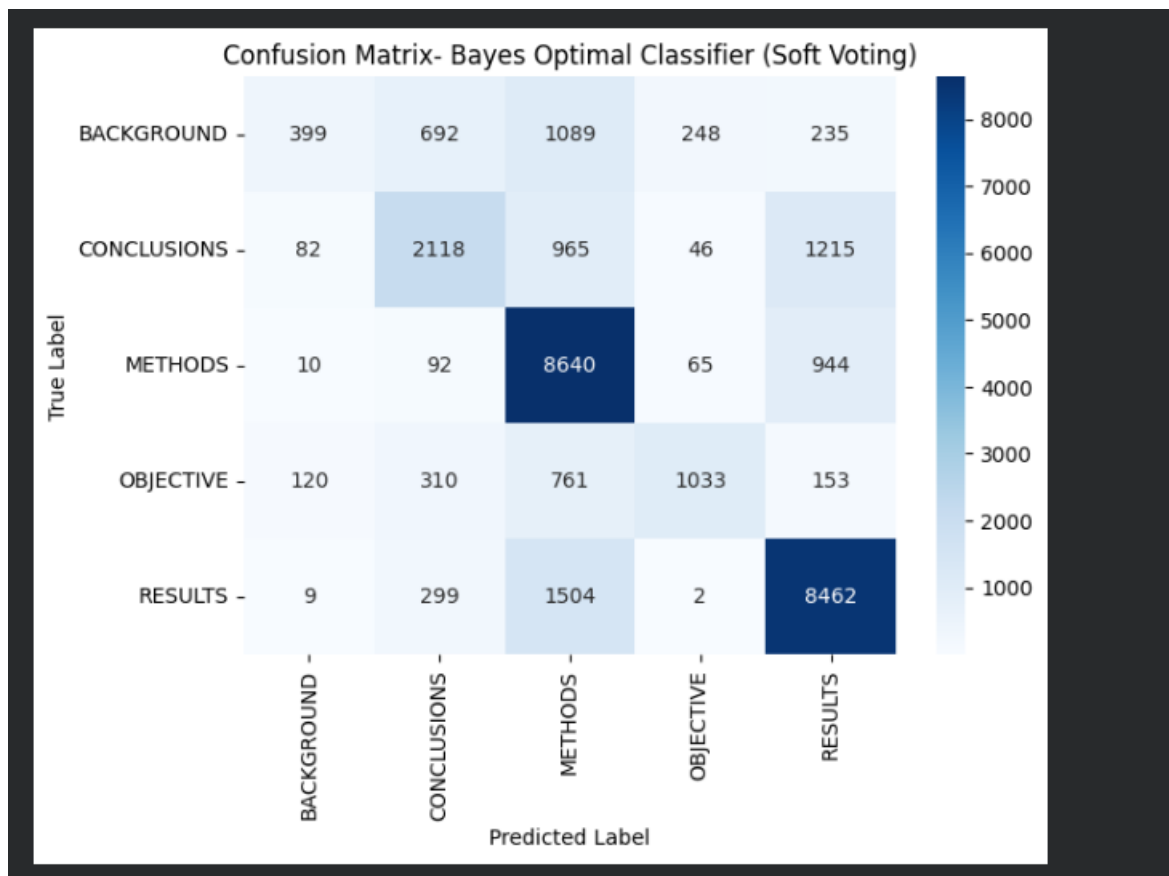
=== Final Evaluation: Bayes Optimal Classifier (Soft Voting) ===

BOC Accuracy: 0.7002

BOC Macro F1 Score: 0.5761

Classification Report:

	precision	recall	f1-score	support
BACKGROUND	0.64	0.15	0.24	2663
CONCLUSIONS	0.60	0.48	0.53	4426
METHODS	0.67	0.89	0.76	9751
OBJECTIVE	0.74	0.43	0.55	2377
RESULTS	0.77	0.82	0.80	10276
accuracy			0.70	29493
macro avg	0.68	0.55	0.58	29493
weighted avg	0.70	0.70	0.67	29493



The comparison shows that while individual Naive Bayes models provide efficient and interpretable results, the soft voting ensemble delivers superior overall performance. By combining predictions from multiple base models and weighting them by their posterior probabilities, the soft voting approach effectively reduces individual model biases and improves generalization. This demonstrates that ensemble learning can significantly enhance classification accuracy and robustness compared to using a single model alone.