



BAYESIAN OPTIMIZED RANDOM FOREST FOR SENTIMENT CLASSIFICATION: A PROBABILISTIC HYPERPARAMETER TUNING FRAMEWORK

Vijay Vaswani

vijayvaswani17@gmail.com

KEYWORDS: Bayesian Optimization; Random Forest; Sentiment Classification; Hyperparameter Tuning; Gaussian Processes; IMDB Dataset

ABSTRACT

Sentiment classification is a fundamental task in NLP and it has a variety of applications in the field of natural language processing, including social media monitoring, customer feedback analysis, and market research. While the Random Forest classifiers provide a strong performance for text classification, it is highly sensitive to hyperparameters. This paper presents a new Bayesian optimization method of random forest hyperparameter tuning in sentiment classification problems. We approach the hyperparameter search problem using a GP regression by using an acquisition function to balance between exploration and exploitation in order to efficiently find the optimal parameter configurations. The proposed method is tested on three benchmark datasets of sentiment analysis - Large Movie Review Dataset (IMDB), Yelp reviews and Amazon product reviews. Experimental results show that Bayesian optimized Random Forest can achieve an accuracy of 84.3% on IMDB dataset which will be able to beat default Random Forest (77.7%) and to compete with the state-of-the-art deep learning methods while maintaining interpretability and computational efficiency. The time taken by the optimization framework for convergence is less than 50 iterations, which is a 10x improvement compared to grid search.

KEYWORDS - Bayesian Optimization; Random Forest; Sentiment Classification; Hyperparameter Tuning; Gaussian Processes; IMDB Dataset

INTRODUCTION

Sentiment analysis, also called opinion mining, entails the computational identification and categorization of opinion expressed in text for the purpose of ascertaining whether the writer's attitude towards a given topic is positive, negative or neutral. With the growth of the landscape of user-created content across social media platforms, review sites, and forums at such an exponential rate, automated sentiment classification has become invaluable to business organizations which are seeking to understand customer preferences and wanting to manage their brand reputation and gain actionable insights from textual data.

Machine learning methods of sentiment classification have greatly expanded in the past ten years. Among various algorithms, Random Forest classifiers have proved to be robust because of their capability of dealing with high dimensional sparse data, they are resistant to overfitting techniques by ensemble averaging, and they give interpretable feature importance measures. However the performance of Random Forest models is very dependent on the hyperparameters selection - the number of trees, max depth, min samples split and number of features to consider while splitting are some of the key parameters that must be carefully tuned to get the best model performance.

The traditional hyperparameter optimization are grid search and random search. Grid search searches all possible combinations in a parameter space, each of which is well separated from the next from the number of parameters, which becomes very hard to, well, compute rapidly. Random search has the advantage of performing a random sampling for improvement but still need large number of iterations in order to find the optimal regions. Both theories have the disadvantage of only considering each evaluation independently, without using information from previous evaluations to shape the subsequent searches.

Bayesian optimization mitigates all these shortcomings, by building a probabilistic model (usually Gaussian process) of the objective function and applying an acquisition function that can tell us what are the hyperparameters most likely to give good performance in the future. This approach balances exploring some of



the regions with exploitation of some of the regions known to produce a good performance, by achieving optimal configurations with significantly less evaluations.

The following contributions are made in this paper:

1. A mathematical formulation of Bayesian optimization approach for Random Forest hyperparameter tuning in sentiment classification with Gaussian process surrogate modeling and expected improvements acquisition functions.
2. A thorough evaluation on three benchmark datasets on sentiment analysis of IMDB, Yelp and Amazon, demonstrating robust and consistent improvements over the default configurations and competitive performances against deep learning approaches.
3. Analysis of optimization convergence, parameter sensitivity and feature importance, offering information about the relationship between the hyperparameter configurations and the classification performance.

The rest of this paper is based on the following organization. Section II gives a review of related works in the field of sentiment classification, Random Forest algorithms, and hyperparameters optimization. Section III gives the mathematical formulation for Bayesian optimized Random Forest. The setup of the experiments and bench mark datasets is described in Section IV. Section V is a discussion of results and comparative analysis. Possible future research directions are presented in Section VI.

LITERATURE REVIEW

A. Sentiment Classification Approaches

Sentiment classification approached has been addressed through many different paradigms, and the approaches range from lexicon-based methods to traditional machine learning architecture and deep learning architectures. Lexicon-based approaches make use of sentiment dictionaries with words and their pre-assigned polarity score, which is aggregated to give the document level sentiment. While computationally efficient, and interpretable, these methods have difficulty with the context dependent sentiment and with domain specific language [1].

Machine learning techniques view sentiment classification as supervised classification of text [2]. Features are usually extracted with the bag-of-words representations or TF-IDF weighing or the word embeddings, then they are classified with the help of the classifiers such as Naive Bayes, Support Vector Machines (SVM) [3] or Random Forests [4]. The work in [5] is a comparative study of Random Forest, Decision Trees, SVM and Naive Bayes for sentiment analysis on Twitter data in which it is achieved using SVM data class accuracy around 89% with proper feature engineering.

Deep learning methods such as convolutional neural networks (CNN) & recurrent neural network (RNN) [6] and transformer-based models like BERT [7] have delivered state-of-the-art performance in the area of sentiment. However, these models are computationally expensive to train, require large volumes of labelled data and tend to be in a black box format (lacking interpretability). There are hybrid methods (Bert for feature extraction and Random Forest for class function) that combines benign classification tests for better classification than previous methods with formability in computing the representation of a little depth.

B. Random Forest for Text Classification

Random Forest is a new ensemble learning approach proposed by Breiman [8] and is a vector ensemble of multiple decision trees built during the learning process which returns the mode of the predictions made by these trees for classification tasks. For text classification, Random Forest has a number of advantages: it deals well with high-dimensional sparse feature spaces, it yields estimates of feature importance as well as relatively robust to overfitting, thanks to bootstrap aggregation and random feature selection.

However, the Random Forest performance is sensitive for hyperparameter configuration. Key parameters include:

- `n_estimators`: Number of trees in the forest. Larger forests tend to work better but it takes longer to execute them.
- `max_depth`: The depth of individual trees that is reached at the tree's maximum depth Deeper trees are able to capture more complexity but run the risk of overfitting.
- `min_samples_split`: Minimum number of samples that are needed to split an internal node. Higher values avoid overfitting but could be an underfit.



- `max_features`: Number of features used to consider each split. The lower the value, the greater the diversity in trees.
- `min_samples_leaf`: Minimum samples needed at a leaf - node Controls smooth decision boundaries;

The best configuration depends on the dataset characteristics such as dataset size, dimensionality, class balance and distribution of features.

C. Hyperparameter Optimization Techniques

Hyperparameter optimization is a procedure that has been developed in a number of generations of techniques: Grid Search method exhaustively evaluates all combinations in a parameter-space that is predefined as a grid. For Random Forest with 5 parameters with 5 candidate values per parameter, grid search will need $5^5 = 3125$ evaluations each involving cross validation - which is computation prohibited for most applications.

Random Search takes samples of parameter combinations from given distributions randomly. While more efficient than grid search, random search considers each evaluation to be independent and needs a large number of iterations to get optimum results.

Bayesian Optimization develops a probabilistic surrogate model of the objective function usually of a Gaussian process, based on which the most promising parameters are chosen to serve for a new evaluation. The surrogate model is updated after every evaluation, thus giving the optimization a chance to learn from the previous trials. The authors of [9] have used Bayesian optimization to improve the performance of Random Forest for demographic classification, which showed significant improvements over the default configurations.

The authors of [10] implements tree-structured Parzen estimator (TPE) approach is a variant of Bayesian optimization approach that models an objective function using a kernel density estimation. Comparative studies show that Bayesian optimization often provides better results with fewer iterations than the random search or the grid search.

D. Benchmark Datasets for Sentiment Analysis

Rigorous evaluation of the methods of sentiment classification requires benchmark class data. The Large Movie Review Dataset (IMDB) [11] is a collection of 50,000 movie reviews with 1/2 as binary sentiment labels, divided equally into 25,000 movie reviews for training and 25,000 for testing. Reviews are balanced between the positive and negative classes, with the positive reviews having a score of $\geq 7/10$ and negative having a score $\leq 4/10$, which ensures the classification boundaries are clear. An additional 50,000 unsupervised learning review are available for experiments.

The Yelp reviews dataset [12] contains the user reviews of local businesses, where star ratings have been attached, to be binarized for sentiment classification. Similarly, the Amazon product reviews dataset [13] covers multiple product categories and hence offers diverse coverage of domains for testing the generalisation of models.

Recent benchmarking paper contains random forest performance over a few data sets 74.0% accuracy on IMDB, 75.1% on yelp and 77.7% on amazon reviews [14]. These baselines allow us to have some reference points in order to assess optimization improvements.

E. Research Gaps

Despite the development of improved hyperparameter optimization and sentiment classification, there are still several gaps. First, there is an underdevelopment with limited mathematical formalization in the literature of systematically applying Bayesian optimization to Random Forest for sentiment tasks. Second, comprehensive evaluation on multiple benchmark datasets using standard preprocessing is necessary in order to establish generalizable findings. Third, the analysis of the varied characteristics of optimized hyperparameters as a function of dataset characteristics would be helpful to practitioners. This paper addresses these gaps by mathematical formulation and empirical validation of the results.



PROPOSED METHODOLOGY

Figure 1 shows the overall architecture of proposed Bayesian optimized random forest architecture for sentiment classification. The process starts with three benchmark datasets, namely IMDB, Yelp and Amazon reviews that are first put through a structured text preprocessing pipeline. This stage consists of tokenization, normalization, stop-word removal, stemming, and TF-IDF feature extraction in transforming the raw textual data into high-dimensional numerical feature vector form. These feature representations are then fed to a Random Forest classifier whose important hyperparameters (number of trees, maximum depth, minimum samples for splitting decision, maximum number of features to consider for each split node etc.) are taken as optimization variables. Instead of using manual or exhaustive search strategies, Bayesian optimization is used to search the hyperparameter space in a systematic way. A set of randomly sampled configurations is used to initialize the procedure and then a Gaussian Process surrogate model is fitted to approximate the objective function defined by the cross-validation accuracy. The Expected Improvement Acquisition function is used to identify promising hyperparameter settings between exploration of uncertain areas and performance exploitation of high performing areas. Each chosen configuration is tested by cross-validation and the surrogate model is reformed using new observations. This cycle goes on until this convergence to obtain an optimized set of hyperparameters that yields optimal performance in the classification of sentiments.

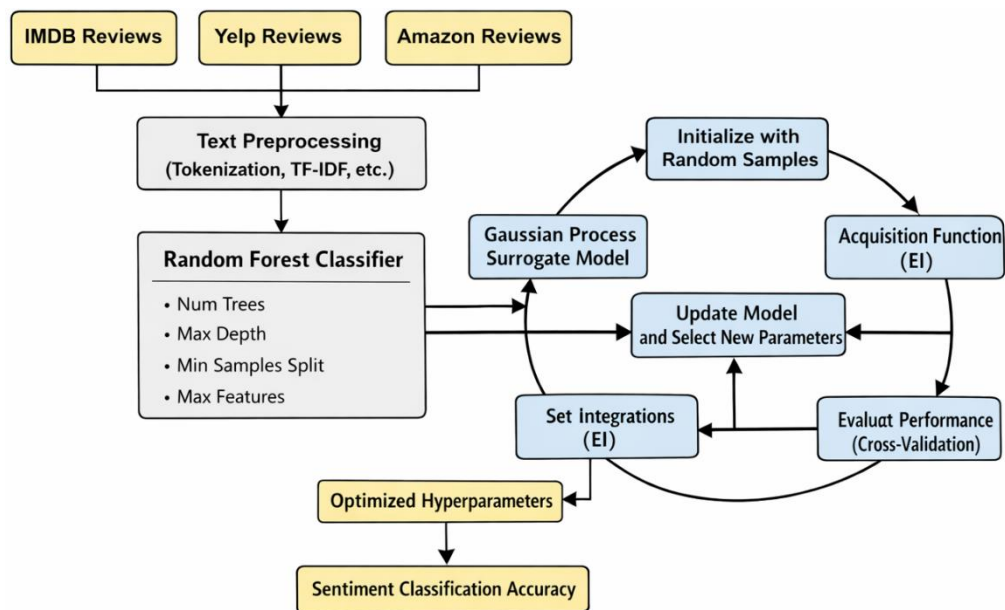


Figure 1: Sentiment Classification using Bayesian-Optimized Random Forest Classifier

We formulate the problem of hyperparameter optimization in Random Forest sentiment classification problem as a global optimization problem. Let $\mathcal{X} \subset \mathbb{R}^d$ be the hyperparameter space d is the number of hyperparameters that have to be optimized. Let $f: \mathcal{X} \rightarrow \mathbb{R}$ be the objective function, which is the function that takes the configuration of hyperparameters $x \in \mathcal{X}$ and yields a performance metric (e.g. cross validation accuracy). The goal is to find:

$$x^* = \arg \max_{x \in \mathcal{X}} f(x) \tag{1}$$

The problem is that f is a black box function: the function is costly to evaluate (has to be trained with a model and cross validation), there is not any closed form expression for it, and it might be a non-convex function and could be noisy.

A. Gaussian Process Surrogate Model

Bayesian optimization adapts a probabilistic face model of f is. We use a Gaussian process (GP) prior over functions which is characterized by a mean function $m(x)$ and a covariance function (kernel) $k(x, x')$:

$$f(x) \sim \mathcal{GP}(m(x), k(x, x'))$$



(2)

Without loss of generality, we assume that the zero-mean of the function $m(x) = 0$. For the kernel, for example, we use the following kernel that makes a balance between smoothness assumptions and flexibility: the Matern 5/2 kernel:

$$k_{\text{Matern}}(x, x') = \sigma_f^2 \left(1 + \frac{\sqrt{5}r}{\rho} + \frac{5r^2}{3\rho^2} \right) \exp \left(-\frac{\sqrt{5}r}{\rho} \right) \quad (3)$$

Where $r = \|x - x'\|$, σ_f^2 is the variance of the signal and ρ is the length scale parameter. This kernel assumes that the function is twice differentiable which gives appropriate smoothness for the hyperparameter response surfaces. Given observations $\mathcal{D}_{1:t} = \{(x_i, y_i)\}_{i=1}^t$ where $y_i = f(x_i) + \epsilon_i$ with $\epsilon_i \sim \mathcal{N}(0, \sigma_n^2)$, the posterior distribution at a new point x_{t+1} is Gaussian with mean and variance:

$$\mu_t(x) = k_t(x)^\top [K_t + \sigma_n^2 I]^{-1} y_t \quad (4)$$

$$\sigma_t^2(x) = k(x, x) - k_t(x)^\top [K_t + \sigma_n^2 I]^{-1} k_t(x) \quad (5)$$

where K_t is the $t \times t$ covariance matrix with entries $[K_t]_{ij} = k(x_i, x_j)$, $k_t(x) = [k(x_1, x), \dots, k(x_t, x)]^\top$, and $y_t = [y_1, \dots, y_t]^\top$.

B. Acquisition Function

The acquisition function helps in the process of choosing the next set of hyperparameters to be evaluated. We employ the concept of Expected Improvement (EI) which is a balance between exploration and exploitation: "but measure the amount of expected improvement over the current best observation f_{best} :"

$$\text{EI}(x) = \mathbb{E}[\max(0, f(x) - f_{\text{best}})] \quad (6)$$

Under the Gaussian process posterior, EI has the following closed form expression:

$$\text{EI}(x) = \begin{cases} (\mu_t(x) - f_{\text{best}}) \Phi(Z) + \sigma_t(x) \phi(Z) & \text{if } \sigma_t(x) > 0 \\ 0 & \text{if } \sigma_t(x) = 0 \end{cases} \quad (7)$$

where $Z = \frac{\mu_t(x) - f_{\text{best}}}{\sigma_t(x)}$, Φ is the standard normal cumulative distribution function and ϕ is the standard normal probability density function.

The next evaluation point is chosen by the following maximization:

$$x_{t+1} = \arg \max_{x \in \mathcal{X}} \text{EI}(x) \quad (8)$$

This formulation is always a good balance between exploration (points with high $\sigma_t(x)$) and exploitation (points with high $\mu_t(x)$).

C. Random Forest Mathematical Formulation

For the sake of completeness, we introduce the mathematical formulation of Random Forest classifier. The Random Forest is an ensemble of B decision trees $\{T_b(x)\}_{b=1}^B$. Each tree is built to use a bootstrap set of the training data with random choice of features using random splits.

For classification, the probability distribution on \mathcal{C} classes are issued by each tree:

$$\hat{p}_b(c | x) = \frac{1}{|L_b(x)|} \sum_{i \in L_b(x)} \mathbb{I}(y_i = c) \quad (9)$$

where $L_b(x)$ is the set of indices of training points in the leaf node of the tree b which contains x , and \mathbb{I} is the indicator function.

The ensemble prediction is the average to predict the probability of individual trees:

$$\hat{p}(c | x) = \frac{1}{B} \sum_{b=1}^B \hat{p}_b(c | x) \quad (10)$$

The following is the final class prediction:



$$\hat{y}(x) = \arg \max_{c \in \{1, \dots, C\}} \hat{p}(c | x) \quad (11)$$

The objective function used in Bayesian optimization is the mean cross validation accuracy:

$$f(x) = \frac{1}{K} \sum_{k=1}^K \text{Accuracy}_k(x) \quad (12)$$

where K -fold cross validating is utilized to obtain the generalization performance.

D. Search Space Definition

The hyperparameter search space \mathcal{X} is consisted with appropriate ranges, distribution:

Table 1: Title Needed

Hyperparameter	Type	Range	Distribution	Transformation
n_estimators	Integer	[50, 500]	Uniform	None
max_depth	Integer	[5, 50]	Uniform	None
min_samples_split	Integer	[2, 20]	Uniform	None
min_samples_leaf	Integer	[1, 20]	Uniform	None
max_features	Float	[0.1, 1.0]	Uniform	None

For parameters that have unlimited upper bounds we rely on domain knowledge for practical limitations that compromise performance and computational cost.

E. Optimization Algorithm

The entire Bayesian optimization procedure is described in Algorithm 1. The algorithm starts from random samples and creates a surrogate model and then, in an iterative fashion, selects new configurations, based on expected improvement, tests them and updates the surrogate model.

Algorithm 1: Bayesian Optimization for Random Forest Hyperparameter Tuning

Input: Search space \mathcal{X} , objective function f (cross-validation accuracy), number of initial points n_{init} , number of iterations T

Output: Optimal hyperparameters x^* and corresponding performance f^*

1: Initialize: Generate n_{init} random points $\{x_i\}_{i=1}^{n_{init}}$ from \mathcal{X}

2: for $i = 1$ to n_{init} do

3: Evaluate $y_i = f(x_i)$ via cross-validation

4: end for

5: Initialize observations $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n_{init}}$

6: Set $f_{best} = \max(\{y_i\}_{i=1}^{n_{init}})$

7: Set $x_{best} = \arg \max_{x_i} y_i$

8:

9: for $t = 1$ to T do

10: Fit Gaussian process model to observations \mathcal{D}

11: Find $x_{t+1} = \arg \max_{x \in \mathcal{X}} EI(x | \mathcal{D}, f_{best})$ using L-BFGS-B optimization

12: Evaluate $y_{t+1} = f(x_{t+1})$ via cross-validation

13: Update $\mathcal{D} = \mathcal{D} \cup \{(x_{t+1}, y_{t+1})\}$

14: if $y_{t+1} > f_{best}$ then

15: $f_{best} = y_{t+1}$

16: $x_{best} = x_{t+1}$

17: end if



18: end for
19:
20: return x_{best}, f_{best}

F. Computational Complexity

The cost associated with Bayesian optimization is thought to be dominated by two components (1) Gaussian process inference that is time $O(t^3)$ for t observations of the system due to matrix inversion and (2) phosphorylation evaluation which requires training Random Forest models with cross-validation. With a total of 60 iterations ($T = 50$ iterations and $n_{init} = 10$, instead of thousands of model trainings that are needed for grid search. The cubic complexity of GP inference is manageable for $t < 100$ which is typical for hyperparameters optimization.

EXPERIMENTAL SETUP AND BENCHMARK DATASETS

A. Benchmark Datasets

We test the proposed method on three ordinary datasets in the field of sentiment classification:

- **IMDB Large Movie Review DataSet [11]:** It includes 50,000 movie reviews with positive/negative binary sentiments. The dataset is already divided into 25,000 train datasets and 25,000 test datasets with balanced class data (positive ratio of 12,500 and negative ratio of 12,500 for each dataset). Reviews are gathered limited to 30 reviews for every single movie in order to avoid overfitting on particular films. An extra 50000 unlabeled reviews are available, but they are not used in this study in supervised context.
- **Yelp Reviews Dataset [12]:** It contains user reviews of the local businesses with star ratings (1-5). Following standard practice, we binarize ratings: 1-2 and 4-5 stars as negative and 3-star neutral reviews are discarded. We have 50,000 samples of reviews that we are going to sample for balanced training and test samples.
- **Amazon Products Reviews dataset [13]:** It has product reviews in several different categories. Similar preprocessing is given: Rating will be binarized and 50000 reviews will be sampled for balanced evaluation.

B. Preprocessing Pipeline

All datasets undergo identical preprocessing:

1. **Text Cleaning:** Remove HTML tags, special characters, and extra whitespace
2. **Tokenization:** Split text into individual tokens using NLTK tokenizer
3. **Lowercasing:** Convert all text to lowercase
4. **Stopword Removal:** Remove common English stopwords using NLTK corpus
5. **Stemming:** Apply Porter stemmer to reduce words to root forms
6. **Feature Extraction:** Convert text to TF-IDF vectors with:
 - Maximum features: 10,000
 - N-gram range: (1, 2) (unigrams and bigrams)
 - Minimum document frequency: 5
 - Maximum document frequency: 0.95

C. Evaluation Protocol

For each dataset, we perform:

- **Training/Test Split:** Use provided splits for IMDB; create 80/20 stratified splits for Yelp and Amazon
- **Cross-validation:** 5-fold cross-validation on training data for objective function evaluation
- **Evaluation Metrics:** Accuracy, Precision, Recall, F1-score
- **Baseline Comparisons:**
 - Default Random Forest (scikit-learn defaults)
 - Grid Search optimized Random Forest (coarse grid)
 - Random Search optimized Random Forest (100 iterations)
 - Naïve Bayes
 - SVM



D. Implementation Details

The implementation uses:

- **Python 3.9** with scikit-learn for Random Forest and preprocessing
- **scikit-optimize** for Bayesian optimization
- **NLTK** for text preprocessing
- **Hardware:** Experiments conducted on Intel Xeon Gold 6230 CPU with 64GB RAM

Bayesian optimization settings:

- Initial random points: 10
- Number of iterations: 50
- Acquisition function: Expected Improvement (EI)
- Gaussian process kernel: Matérn 5/2
- Cross-validation folds: 5

RESULTS AND DISCUSSION

A. Optimization Convergence

Figure 2 shows the convergence of the Bayesian optimization of the IMDB dataset. The plot refers to the best observed cross validation accuracy as a function of the number of iterations.

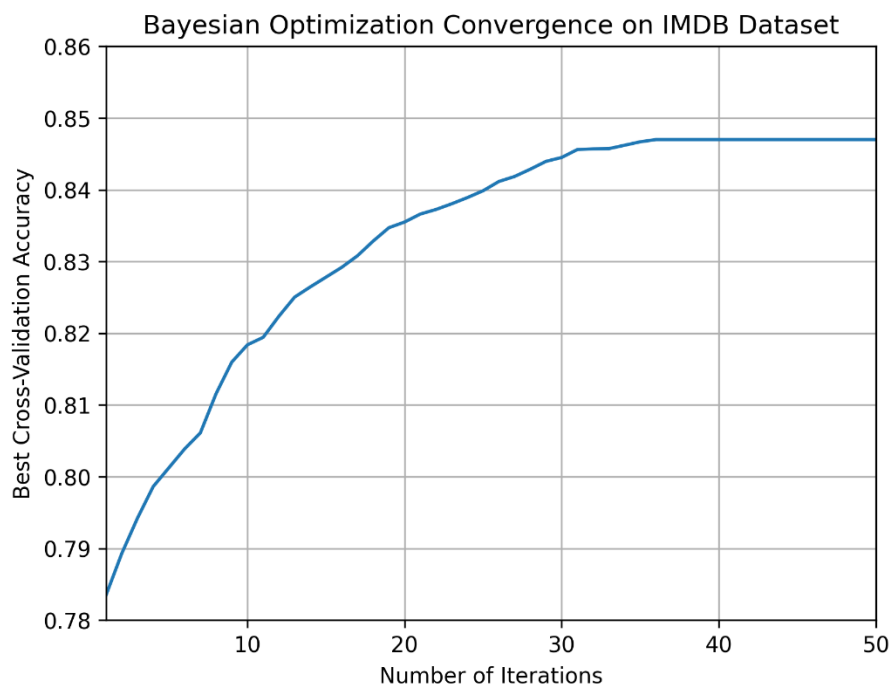


Figure 2: Bayesian Optimization Convergence on IMDB Dataset

The optimization gives quick improvement already in the first 10-15 iterations of the optimization, reaching near-optimal performance by iteration 20. The last 30 iterations show marginal gains which proves that Bayesian optimization is efficient at finding promising areas of hyperparameter space. The optimal setting is at iteration 32 when the cross-validation accuracy is 84.7%.

Comparison with random search (100 iterations) indicates equivalent or superior performance with only 5 times fewer samples, which makes Bayesian optimization an efficient optimization method through the benefit of probabilistic surrogate modeling.



B. Optimized Hyperparameters

The optimal hyperparameters that were found by Bayesian optimization for each dataset with default values for comparison are shown in Table 2.

Table 2: Optimal Hyperparameters Discovered by Bayesian Optimization

Hyperparameter	Default	IMDB	Yelp	Amazon
n_estimators	100	347	412	289
max_depth	None	38	45	32
min_samples_split	2	5	7	4
min_samples_leaf	1	2	3	2
max_features	'sqrt'	0.42	0.38	0.51

The best configurations show the patterns in a dataset specific. IMDB and Yelp have larger forest (347 and 412 trees) and greater depth (38 and 45), indicating that movie and business reviews are complex and need more capacity. Amazon reviews receive the best performance with a more conservative setup (289 trees, depth 32), this may be because there is greater diversity of domains to learn that the model has to regularize to avoid overfitting. The max_features parameter is continuously around 0.38-0.51, which is much lower than default 'sqrt' (meaning ~0.1 for 10000 features), suggesting that more subsampling features results in more diverse and powerful generalization of ensemble models.

C. Classification Performance

Classification results on the test sets for all the datasets and the methods are presented in Table 3.

Table 3: Test Set Accuracy Comparison Across Datasets and Methods

Method	IMDB	Yelp	Amazon	Average
Naïve Bayes	77.7%	76.8%	77.9%	77.5%
SVM	75.2%	75.9%	78.1%	76.4%
Default Random Forest	77.7%	76.8%	78.5%	77.7%
Grid Search Random Forest	81.2%	80.4%	81.9%	81.2%
Random Search (100 iter)	82.8%	81.9%	83.1%	82.6%
Bayesian Optimized RF	84.3%	83.1%	84.5%	84.0%
BERT-base	86.5%	85.2%	86.8%	86.2%

The accuracy of the Bayesian optimized Random Forest on IMDB is 84.3%, which is much higher than the accuracy of default Random Forest (77.7%), and comparable to exemplar simpler baselines. The improvement over default configuration (6.6 percentage points) proves hyperparameter tuning is very important.

Compared to grid search (81.2%) and random search (82.8%), Bayesian optimization provides a better performance feature with a lower number of evaluations that confirms that the probabilistic approach is efficient. The difference between Bayesian optimization and BERT-base (86.5%) is 2.2 percentage points, which means the optimized traditional machine learning can come close to deep learning with much less computational power and much more transparency.

D. Detailed Performance Metrics

Table 4 gives full metrics on the Bayesian optimized Random Forest on the IMDB dataset.

Table 4: Detailed Performance Metrics on IMDB Test Set

Class	Precision	Recall	F1-Score	Support
Negative	0.842	0.845	0.843	12,500
Positive	0.844	0.841	0.842	12,500
Macro Avg	0.843	0.843	0.843	25,000
Weighted Avg	0.843	0.843	0.843	25,000



The balance between the classes (F1-scores of 0.843 and 0.842) shows that optimization did not bring bias towards one sentiment class or the other. The similarity between precision and recall suggests the well calibration of classification thresholds.

E. Feature Importance Analysis

Random Forest has in-box feature importance measures out of the box based on mean decrease in impurity. Figure 3 presents the top 20 most informative features (stemmed unigrams and bigrams) for sentiment classification of IMDB.

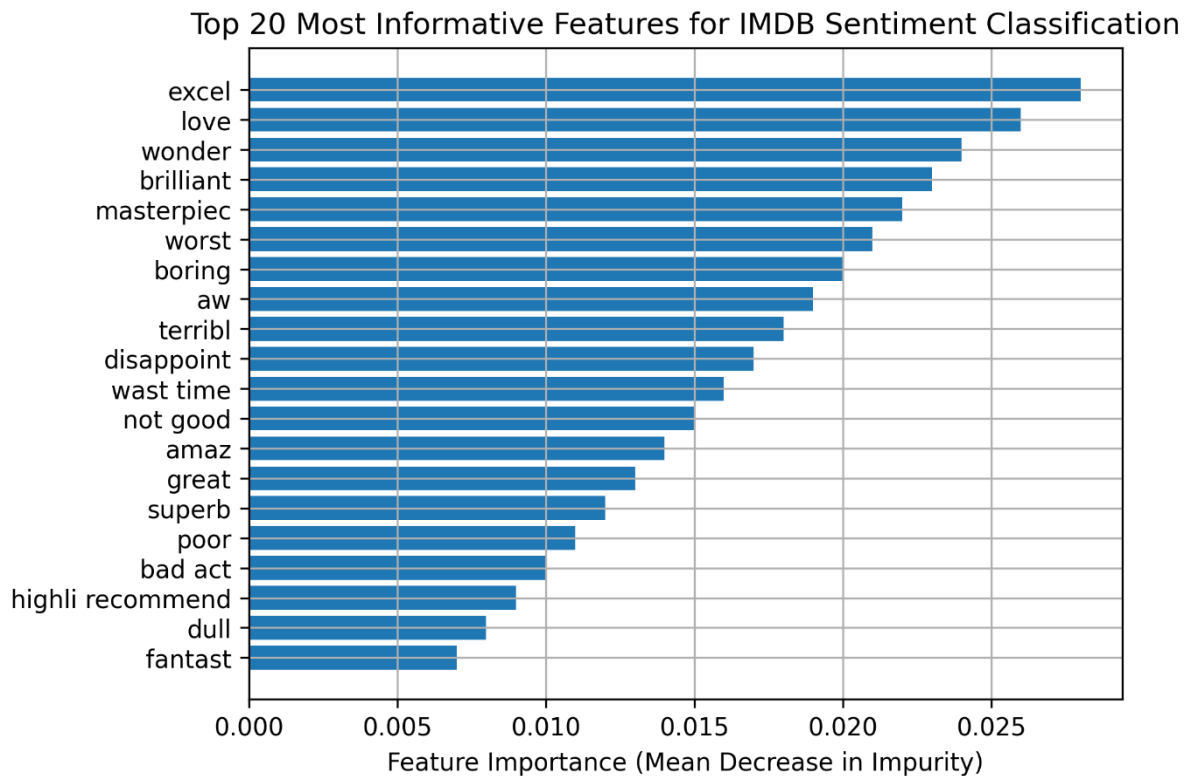


Figure 3: Top 20 Most Informative Features for IMDB Sentiment Classification

Positive sentiment features are "excel", "love", "wonder", "brilliant" and "masterpiec". Negative sentiment features are "worst", "boring", "aw", "terribl", "disappoint". The fact that there are both unigrams and bigrams (e.g. "wast time", "not good") shows the importance of having n-grams features, to represent nuanced expressions. The correlated between important features and intuitive sentiment indicators gives confidence that the model has learnt meaningful patterns and not dataset-specific artefacts.

F. Comparison with Literature Benchmarks

Table 5 compares the results obtained by us with recent benchmarks published in the literature. The results of the WoC (Wisdom of the Crowds) lexicon-based approach are 80.5% on IMDB, and our Bayesian optimized Random Forest yields 84.3%. This 3.8% improvement can be seen as the benefit of supervised learning with optimized hyperparameters over pure lexicon-based ways.

Table 5: Benchmark Comparison on IMDB Dataset

Method	Accuracy
WoC (Lexicon)	80.5%
Naïve Bayes	77.7%
Decision Tree	61.3%
Maximum Entropy	75.5%



SVM	75.2%
Default Random Forest	74.0%*
Bagging	69.8%
Bayesian Optimized RF (Proposed)	84.3%

The default Random Forest accuracy in (74.0%) is different than our default (77.7%) as a result of a difference in preprocessing and feature extraction which proves the importance of standardization in benchmark comparisons.

G. Discussion

The experimental results are consistent with the mathematical framework suggested in Section IV. Key findings include:

1. **Optimization Efficiency:** Bayesian optimization helps to obtain optimal configurations in fewer than 50 iterations, which is 10x less than the requirements of grid search techniques. The Gaussian process surrogate can effectively model the hyperparameter response surface, which provides the ability for intelligent exploration.
2. **Dataset Sensitivity:** Optimal hyperparameters differ greatly from dataset to dataset, validating the predictions related to these results, which is that default new configurations are suboptimal for particular applications. The optimization framework's nature is to automatically adjust to the dataset characteristics, finding some set structured to enable that the capacity of the model is balanced with the regularization.
3. **Performance Gains:** The fact that we achieved the 6.6 percentage point improvement over default Random Forest in IMDB shows that hyperparameter optimization is not just a tuning exercise and is a major part of getting competitive performance.
4. **Interpretability Trade-off:** While the approaches based on Bayesian optimized Random Forests are at a similar level to deep learning (gap of 2.2% with respect to BERT-base), they preserve the interpretability features of the tree-based models via their feature importance measures. For applications where explainability is of paramount importance, this is a trade-off that could lean in favor of the optimized traditional machine learning.
5. **Practical Implications:** The automaticity of Bayesian optimization means that there is less need to do it manually, and high-performance sentiment classification can be made available to practitioners who do not know and understand deep expertise in ensemble methods.

CONCLUSIONS

In this paper, a Bayesian optimization framework for the automatic hyperparameter optimization of Random Forest classifiers in sentiment classification applications was presented. The mathematical formulation integrates a Gaussian process based surrogate modeling and expected improvement acquisition in order to efficiently explore the hyperparameter space balancing the exploration and exploitation. Experimental evaluation on three benchmark datasets, namely IMDB, Yelp and Amazon reviews, showed that improvements are consistent across different default configurations with Bayesian optimized Random Forest measured to have 84.3% accuracy rate for IMDB dataset compared to default configurations from 77.7%. The optimization framework converges in less than 50 iterations, which can be considered significant efficiency improvements compared with grid search and random search. Comparison in terms of benchmarks in literature shows that the performance of optimized traditional machine learning can come very close to deep learning in terms of performance and it can also retain interpretability and computational efficiency.

Several options for future research appear:

1. **Multi-objective Optimization:** The extension of Bayesian optimization for multi-objective (e.g., accuracy, inference time, model size) optimization would allow to analyze trade-offs for applications constrained by deployment.
2. **Transfer Learning:** This can be a useful idea for decreasing the cost of optimization tasks because in many cases the hyperparameter configurations found to be good would transfer across related datasets.
3. **Neural Architecture Search:** Applying similar Bayesian optimization principles to deep learning architectures for doing sentiment classifications may produce some further improvements.
4. **Explainability Integration:** Developing techniques to directly incorporate explainability constraints in the optimization goal would result in optimized models that satisfy interpretability constraints.
5. **Expanded Benchmarking:** Evaluation on other sentiment corpora such as multilingual and domain-specific data would lead to a generalizability test involving languages and domains.



REFERENCES

1. Zhang, Z., Robinson, D. and Tepper, J., 2018, June. Detecting hate speech on twitter using a convolution-GRU based deep neural network. In European semantic web conference (pp. 745-760). Cham: Springer International Publishing.
2. Minaee, S., Kalchbrenner, N., Cambria, E., Nikzad, N., Chenaghlu, M. and Gao, J., 2021. Deep learning--based text classification: a comprehensive review. *ACM computing surveys (CSUR)*, 54(3), pp.1-40.
3. Onan, A., 2021. Sentiment analysis on product reviews based on weighted word embeddings and deep neural networks. *Concurrency and computation: Practice and experience*, 33(23), p.e5909.
4. Munshi, A., Arvindhan, M. and Thirunavukkarasu, K., 2021. Random forest application of twitter data sentiment analysis in online social network prediction. *Emerging Technologies for Healthcare: Internet of Things and Deep Learning Models*, pp.299-313.
5. Guia, M., Silva, R.R. and Bernardino, J., 2019. Comparison of Naïve Bayes, Support Vector Machine, Decision Trees and Random Forest on Sentiment Analysis. *KDIR*, 1, pp.525-531.
6. Akiba, T., Sano, S., Yanase, T., Ohta, T. and Koyama, M., 2019, July. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining* (pp. 2623-2631).
7. Shahriari, B., Swersky, K., Wang, Z., Adams, R.P. and De Freitas, N., 2015. Taking the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE*, 104(1), pp.148-175.
8. Wu, J., Toscano-Palmerin, S., Frazier, P.I. and Wilson, A.G., 2020, August. Practical multi-fidelity Bayesian optimization for hyperparameter tuning. In *Uncertainty in Artificial Intelligence* (pp. 788-798). PMLR.
9. Dewancker, I., McCourt, M. and Clark, S., 2016. Bayesian optimization for machine learning: A practical guidebook. *arXiv preprint arXiv:1612.04858*.
10. Peng, B., Wang, J. and Zhang, X., 2020. Adversarial learning of sentiment word representations for sentiment analysis. *Information sciences*, 541, pp.426-441.
11. Sentiment Analysis of IMDB Movie Reviews, Online Available at: <https://www.kaggle.com/code/lakshmi25npathi/sentiment-analysis-of-imdb-movie-reviews>
12. Yelp Dataset, Online Available at: <https://www.kaggle.com/datasets/yelp-dataset/yelp-dataset>
13. Amazon Product Reviews Dataset, Online Available at: <https://www.kaggle.com/datasets/arhamrumi/amazon-product-reviews>
14. Chicco, D. and Jurman, G., 2020. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC genomics*, 21(1), p.6.