

Vietnam National University - Ho Chi Minh City
Ho Chi Minh City University of Technology
Faculty of Computer Science and Engineering



MACHINE LEARNING

Assignment Report

COMPREHENSIVE MACHINE LEARNING MODELS: IMPLEMENTATION AND COMPARISON

Instructor: Nguyễn An Khương, CSE-HCMUT

Student(s): Lê Quốc Bảo - 2252065 - Leader
Lê Đỗ Minh Anh - 2252023
Nguyễn Quốc Anh - 2252035

Ho Chi Minh City, 3/2025



Contribution

No.	Name	Student ID	Task	Contribution
1	Lê Quốc Bảo	2252065	Theory, implementation, and detailed explanation for Decision Tree and Neural Network. Writing overall report.	33%
2	Lê Đỗ Minh Anh	2252023	Theory, implementation, and detailed explanation for Graphical Models (Bayesian Networks, HMM). Running and evaluating model results.	33%
3	Nguyễn Quốc Anh	2252035	Theory, implementation, and detailed explanation for Naive Bayes. Creating and editing Github Repository.	33%

Table 1: Project Contributions

Contents

1	Overview	3
2	Data Preprocessing	3
2.1	Loading the Dataset	3
2.2	Mapping Labels	3
2.3	Handling Missing Values	3
3	Pipeline Construction	4
3.1	Text Preprocessing	4
3.2	Model Implementation	4
3.2.1	Decision Trees	4
3.2.2	Naïve Bayes	4
3.2.3	Neural Networks	4
3.2.4	Graphical Models	5
3.3	Pipeline Integration	5
4	Hyperparameter Tuning	5
5	Performance Analysis	6
5.1	Evaluation Metrics	6
5.2	Cross-Validation	6
6	Model Analysis	6
6.1	Decision Tree Model	6
6.1.1	Overview	6
6.1.2	Key Parameters and Their Impact	7
6.1.3	Feature Importance Analysis	7
6.1.4	Strengths and Limitations	8
6.2	Neural Network	8
6.2.1	Overview	8

6.2.2	Key Parameters and Their Impact	9
6.2.2.a	MLP Classifier Parameters	9
6.2.2.b	TF-IDF Vectorizer Parameters	9
6.2.3	Feature Importance Analysis	9
6.3	Strengths and Limitations	10
6.3.1	Strengths	10
6.3.2	Limitations	10
6.4	Naive Bayes	11
6.4.1	Overview	11
6.4.2	Key Parameters and Their Impact	11
6.4.3	Feature Importance Analysis	11
6.4.4	Strengths and Limitations	12
6.5	Graphical Model	12
6.5.1	Overview	12
6.5.2	Key Parameters and Their Impact	13
6.5.2.a	Bayesian Network	13
6.5.2.b	Hidden Markov Model	13
6.5.3	Feature Importance Analysis	14
6.5.3.a	Bayesian Network	14
6.5.3.b	Hidden Markov Model	15
6.5.4	Strengths and Limitations	15
6.5.4.a	Bayesian Network	15
6.5.4.b	Hidden Markov Model	15
7	Usage Example	16
8	Experimental Results	16
8.1	Evaluation Method	16
8.2	Experimental Results and Discussion	17
8.2.1	Neural Network (MLP)	17
8.2.2	Naive Bayes	17
8.2.3	Bayesian Network	18
8.2.4	Decision Tree	18
8.2.5	Graphical Models	18
8.2.6	Overall Discussion	18
9	Conclusion	19
10	References	19

1 Overview

In this project, we evaluate a variety of machine learning models introduced in the course to address a challenging text classification task. The primary objective is to classify news articles into four distinct categories—**World**, **Sports**, **Business**, and **Science/Technology**—using the **AG News** dataset.

The **AG News** dataset serves as a robust benchmark for text classification, comprising thousands of news articles sourced from reputable outlets such as **Reuters** and **AP**. Each article is pre-labeled into one of the four categories, providing a clear and organized division of topics. The dataset is partitioned into training and testing sets, ensuring a balanced and realistic scenario for evaluating model performance. Moreover, the concise headlines and summaries further enhance its utility for developing and comparing various text classification approaches.

Multiple models were implemented and rigorously tested to analyze their effectiveness in this use case. Each model underwent a comprehensive pipeline, including data preprocessing, hyperparameter tuning, and performance evaluation using standard metrics such as **Accuracy**, **Precision**, **Recall** and **F1-score**. By systematically comparing these models, we aim to evaluate their respective strengths and limitations, ultimately guiding improved model selection and optimization strategies for text classification tasks.

2 Data Preprocessing

2.1 Loading the Dataset

The dataset is loaded from Hugging Face using the `load_dataset` function. Both training and testing sets are converted to pandas DataFrames.

```
1 from datasets import load_dataset
2 dataset_name="fancyzhx/ag_news"
3 dataset = load_dataset(dataset_name)
```

2.2 Mapping Labels

Numeric labels are mapped to string categories using a dictionary (`CATEGORY_MAPPING`), ensuring that outputs are human-readable (e.g., 0 → **World**, 1 → **Sports**, etc.).

```
1 CATEGORY_MAPPING = {0: 'World', 1: 'Sports', 2: 'Business', 3: 'Sci/Tech'}
2 train_df['Category'] = train_df['label'].map(CATEGORY_MAPPING)
3 test_df['Category'] = test_df['label'].map(CATEGORY_MAPPING)
```

2.3 Handling Missing Values

In this implementation, missing values are managed by first identifying any null entries in both the text and the label columns. If missing labels are found, the corresponding rows are dropped to prevent inconsistencies in classification. For missing text values, instead of removing rows, we replace them with a placeholder such as *"Unknown"* to maintain the structure of the dataset and avoid losing valuable information. This approach ensures that the feature extraction process, particularly TF-IDF vectorization, remains unaffected by incomplete data. By effectively handling missing values, we enhance the model's ability to generalize well and prevent biases due to data omissions.

3 Pipeline Construction

3.1 Text Preprocessing

```
1 vectorizer = TfidfVectorizer(lowercase=True, stop_words='english', max_features=1000)
2 label_encoder = LabelEncoder()
3
4 X_tfidf = vectorizer.fit_transform(X_train).toarray()
5 y_encoded = label_encoder.fit_transform(y_train)
```

A `TfidfVectorizer` transforms raw text into numerical features. Its configuration includes:

- Lowercasing of text,
- Removal of English stop words,
- Limiting the number of features to 1000 to manage high dimensionality.

3.2 Model Implementation

The model implementation involves selecting a suitable classification algorithm for text classification and adapting it to handle textual features effectively. Different models have unique requirements for optimal performance:

3.2.1 Decision Trees

Decision trees require specific transformations to process high-dimensional text data effectively. In our approach:

- Text data is transformed into numerical features using `TfidfVectorizer`.
- High dimensionality is managed by limiting features to 1000.
- Pruning strategies, such as Cost Complexity Pruning (CCP), are implemented using the `ccp_alpha` parameter to prevent overfitting.

3.2.2 Naïve Bayes

Naïve Bayes classifiers assume feature independence and require handling of sparse features:

- TF-IDF transformation naturally results in sparse features, which suit Naïve Bayes classifiers well.
- Laplace (or additive) smoothing is applied to avoid zero probability issues.
- The model is optimized by tuning the `alpha` parameter, which controls the level of smoothing.

3.2.3 Neural Networks

A simple neural network model is used for text classification:

- Text features are transformed using `TfidfVectorizer`.
- A basic feedforward neural network with a single hidden layer is implemented.
- The model uses an activation function (e.g., ReLU or sigmoid) to introduce non-linearity.
- Optimization techniques such as gradient descent and backpropagation are used for training.

3.2.4 Graphical Models

Graph-based models like Bayesian Networks and Hidden Markov Models (HMMs) require structured probabilistic dependencies:

- Dependencies between words and phrases are modeled using conditional probability tables.
- State transitions for sequences are implemented to capture contextual relationships.
- Structure learning methods are applied to discover patterns in textual data.

3.3 Pipeline Integration

To ensure an end-to-end automated workflow, a machine learning pipeline is constructed, integrating data preprocessing, feature transformation, model training, and evaluation. The pipeline consists of:

- **Text Vectorization:** Raw text is standardized into TF-IDF features for numerical representation.
- **Model Selection:** Based on task requirements, models such as Decision Trees, Naïve Bayes, Neural Networks, or Graphical Models are implemented.
- **Hyperparameter Tuning:** A grid search approach (`GridSearchCV`) is used to identify optimal model parameters.
- **Evaluation Metrics:** The trained model is evaluated using accuracy, precision, recall, F1-score, and confusion matrices.

By integrating all components into a unified pipeline, we streamline the process of training and testing different models efficiently, enabling fair comparisons and model selection for the best-performing.

4 Hyperparameter Tuning

We demonstrate hyperparameter tuning on Decision Tree model.

```
1 def tune_hyperparameters(pipeline, X_train, y_train):
2     param_grid = {
3         "clf__max_depth": [10, 20, None],
4         "clf__min_samples_split": [2, 5, 10],
5         "clf__ccp_alpha": [0.0, 0.001, 0.01]
6     }
7     skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
8     grid_search = GridSearchCV(pipeline, param_grid, cv=skf, scoring='accuracy')
9     grid_search.fit(X_train, y_train)
10    print("Best hyperparameters found:", grid_search.best_params_)
11    return grid_search.best_estimator_
```

Hyperparameter tuning is performed using `GridSearchCV` to optimize model performance. The tuning process involves:

- Defining a parameter grid for `max_depth`, `min_samples_split`, and `ccp_alpha`.
- Utilizing a 5-fold stratified cross-validation (`StratifiedKFold`) to ensure balanced class distributions across splits.
- Training multiple models on different combinations of hyperparameters and selecting the best performing set.
- Analyzing the best hyperparameters, providing insights into the optimal model configuration.

5 Performance Analysis

5.1 Evaluation Metrics

The trained model is evaluated in a separate test set using multiple performance metrics.

```
1 # Evaluate on the test set
2 print("=== Test Set Performance ===")
3 print("Accuracy:", accuracy_score(y_test, y_pred))
4 print("Classification Report:\n", classification_report(y_test, y_pred))
5
6 # Performance Analysis: Plot confusion matrix
7 plot_confusion_matrix(y_test, y_pred, labels=list(CATEGORY_MAPPING.values()))
```

- **Accuracy Score:** Measures the overall accuracy of the predictions.
- **Classification Report:** Provides precision, recall, and F1-score for each category, offering detailed performance insights.
- **Confusion Matrix:** Visualize the highlight misclassification patterns and assess model robustness.

5.2 Cross-Validation

To ensure stability and robustness, cross-validation is performed on the training set:

```
1 # Cross-validation on the training set
2 print("=== Cross-Validation on Training Set ===")
3 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
4 cv_scores = cross_val_score(tuned_pipeline, X_train, y_train, cv=skf, scoring='accuracy'
5                             )
6 print(f"CV Accuracy Scores: {cv_scores}")
7 print(f"Mean CV Accuracy: {np.mean(cv_scores):.4f} +- {np.std(cv_scores):.4f}")
```

- A 5-fold StratifiedKFold cross-validation approach is used to prevent class imbalance.
- CrossValScore is employed to compute accuracy scores across different training splits.
- Mean and standard deviation of cross-validation scores are calculated to assess the model variance.

6 Model Analysis

6.1 Decision Tree Model

6.1.1 Overview

The Decision Tree model is a rule-based, non-parametric learning algorithm that is particularly useful for text classification due to its interpretability and ability to capture nonlinear decision boundaries. The model recursively splits the dataset based on word feature importance, aiming to create homogeneous groups for classification. The key concept behind Decision Trees is Information Gain and Entropy, which help determine the best split at each node.

1. **Entropy:** Measures the impurity of a dataset:

$$H(S) = - \sum_{i=1}^e P_i \log_2 P_i$$

where P_i is the proportion of class i in the dataset S .

2. **Information Gain:** Measures the reduction in entropy after splitting on a feature:

$$IG(S, A) = H(S) - \sum_{v \in A} \frac{|S_v|}{|S|} H(S_v)$$

- A is the feature being split on.
- S_v is the subset of S where feature A has value v

For this task, the model is trained using a TF-IDF feature representation, which transforms raw text data into numerical values based on term frequency and inverse document frequency. This helps the model identify important words while reducing the impact of common, less informative terms.

6.1.2 Key Parameters and Their Impact

The Decision Tree classifier includes several hyperparameters that directly affect its performance, generalization, and interpretability:

```
1 DecisionTreeClassifier(max_depth=max_depth,
2                         min_samples_split=min_samples_split,
3                         ccp_alpha=ccp_alpha,
4                         random_state=42))
```

Max Depth (max_depth)

- Controls the maximum depth of the tree, limiting the number of hierarchical decision splits.
- A deeper tree captures more intricate patterns but risks overfitting.
- A shallower tree generalizes better but may underfit by failing to capture important patterns.

Minimum Samples per Split (min_samples_split)

- Defines the minimum number of samples required to further split a node.
- Lower values create a highly detailed tree but increase sensitivity to noise.
- Higher values enforce more generalized splits, reducing overfitting.

Complexity Parameter for Pruning (ccp_alpha)

- Controls Minimal Cost-Complexity Pruning, which eliminates splits that provide marginal improvement in impurity reduction.
- A higher `ccp_alpha` encourages more aggressive pruning, simplifying the model.
- A lower `ccp_alpha` retains more splits, increasing complexity and potential overfitting.

6.1.3 Feature Importance Analysis

The Decision Tree model allows for feature importance extraction, which helps understand which words contribute most to classification.

- Features with high importance scores correspond to terms that strongly distinguish between categories (e.g., “stocks” for Business, “match” for Sports).
- Words with lower importance may be neutral across categories, contributing less to decision-making.
- By analyzing TF-IDF feature weights, we can identify which words are the strongest predictors for each category and refine feature selection accordingly.

6.1.4 Strengths and Limitations

Strengths

- **Interpretability:** Decision Trees provide clear, human-readable decision rules.
- **Feature Importance Insights:** The model highlights which words are crucial for classification.
- **Non-Parametric:** No assumptions about data distribution are required.

Limitations

- **Overfitting:** A deep tree captures too much detail, leading to poor generalization.
- **Sparse Feature Sensitivity:** Text data transformed via TF-IDF can introduce high dimensionality, making decision boundary learning more complex.
- **Limited Generalization:** Unlike ensemble models (e.g., Random Forest), a single Decision Tree may struggle with borderline classifications.

6.2 Neural Network

6.2.1 Overview

Neural Network is a function approximator that models complex relationships in data using layers of neurons. For text classification, the most common architectures include fully connected networks and recurrent networks (LSTMs, Transformers).

1. Forward Propagation:

- Each neuron computes a weighted sum of inputs, applies an activation function:

$$z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)}$$

$$a^{(l)} = \sigma(z^{(l)})$$

where:

- $W^{(l)}$ and $b^{(l)}$ are the weights and biases of layer l .
- σ is an activation function.
- $a^{(l)}$ is the activation output.

2. Loss Function (Cross-Entropy Loss for Classification):

$$L = - \sum_{i=1}^N y_i \log \hat{y}_i$$

where:

- y_i is the true class label.
- \hat{y}_i is the predicted probability.

3. Backpropagation:

- Gradients are computed using the chain rule:

$$\frac{\partial L}{\partial W^{(l)}} = \frac{\partial L}{\partial a^{(l)}} \cdot \frac{\partial a^{(l)}}{\partial z^{(l)}} \cdot \frac{\partial z^{(l)}}{\partial W^{(l)}}$$

6.2.2 Key Parameters and Their Impact

6.2.2.a MLP Classifier Parameters

- **Hidden Layer Sizes** (`hidden_layer_sizes=(64, 32)`):

The network has two hidden layers with 64 and 32 neurons respectively. Increasing the number of neurons or layers can capture more complex patterns but also raises the risk of overfitting if the dataset is small or noisy. Conversely, a smaller network may underfit if the complexity of the task exceeds the network capacity.

- **Regularization Parameter** (`alpha=0.001`):

This parameter acts as a weight decay to prevent overfitting by penalizing large weights. A smaller alpha may allow more complex patterns (with potential overfitting), while a larger alpha enforces smoother solutions but might underfit.

- **Maximum Iterations** (`max_iter=100`):

The maximum number of epochs for training defines how long the network is allowed to learn. In this case, 100 iterations may be sufficient for a moderately sized dataset, but if the model has not converged, increasing the iterations could improve performance at the cost of additional computation.

- **Learning Rate** (`learning_rate_init=0.01`):

This controls how fast the model adapts to the problem. A high learning rate can lead to rapid learning but might overshoot minima, while a too-low rate might slow down convergence.

- **Batch Size** (`batch_size=32`):

This parameter determines how many samples are processed before updating the model's weights. A balanced batch size can stabilize the training process and optimize performance.

6.2.2.b TF-IDF Vectorizer Parameters

- **Lowercase Conversion and Stop Words**:

Converting text to lowercase and removing common stop words (using `stop_words='english'`) helps reduce noise and focus on the most discriminative terms.

- **Maximum Features** (`max_features=1000`):

Limiting the feature space to the top 1000 words helps reduce dimensionality and computational cost. However, it also means that some potentially useful rare words might be excluded.

Collectively, these parameters influence the model's capacity to learn from text data. A well-tuned combination of these settings is essential to balance model complexity, convergence speed, and generalization ability.

6.2.3 Feature Importance Analysis

Neural networks, such as the MLP used here, do not inherently provide a straightforward interpretation of feature importance compared to tree-based methods. However, several approaches can be considered:

- **TF-IDF Weights**:

Since the text is vectorized using TF-IDF, the vectorizer assigns weights to words based on their frequency and uniqueness. Analyzing these scores can provide insights into which words are considered more important before they enter the MLP.

- **Weight Analysis in the First Layer:**

One might examine the weights connecting the input features (derived from TF-IDF scores) to the first hidden layer. Larger absolute weights may indicate a higher impact of certain features on the network's decisions. However, due to the non-linear transformations and subsequent layers, this analysis is not as straightforward as with linear models.

- **Model-Agnostic Methods:**

Techniques such as permutation importance, SHAP (SHapley Additive exPlanations), or LIME (Local Interpretable Model-agnostic Explanations) can be applied post-hoc to estimate the contribution of individual features. These methods provide a more nuanced view of feature impact despite the black-box nature of neural networks.

In summary, while the MLP classifier does not offer built-in feature importance scores, complementary methods can be used to interpret the role of individual words in the classification process.

6.3 Strengths and Limitations

6.3.1 Strengths

- **Simplicity and Efficiency:**

The pipeline leverages a straightforward TF-IDF vectorization followed by an MLP classifier. This simplicity allows for efficient training and easy modifications.

- **Nonlinear Learning Capability:**

The neural network can model complex, non-linear relationships in the data, which may result in better performance compared to linear models for certain text classification tasks.

- **Multi-Class Handling:**

The design naturally supports multi-class classification, making it suitable for datasets like AG News where multiple topics are present.

6.3.2 Limitations

- **Limited Interpretability:**

Unlike models such as decision trees, MLPs are less transparent. The lack of direct feature importance scores means that additional analysis (using permutation importance, SHAP, or LIME) is necessary to interpret the model's decisions.

- **Hyperparameter Sensitivity:**

The performance of the model is heavily dependent on the selection of hyperparameters (e.g., hidden layer sizes, alpha, learning rate). Suboptimal choices may lead to underfitting or overfitting.

- **Convergence Concerns:**

With a fixed number of iterations (`max_iter=100`), there is a risk that the model might not fully converge, especially if the dataset is large or complex. This could potentially limit the predictive performance.

- **Feature Representation Constraints:**

Relying solely on TF-IDF for feature extraction might not capture semantic context as effectively as modern techniques like word embeddings or transformer-based models. This could affect the quality of the learned representations for nuanced text data.

6.4 Naive Bayes

6.4.1 Overview

Naive Bayes is a probabilistic classifier based on Bayes' Theorem with an assumption of feature independence. It is particularly effective for text classification tasks due to its simplicity and efficiency. The mathematical formulation of Bayes' Theorem is as follows:

$$P(C_k|X) = \frac{P(X|C_k)P(C_k)}{P(X)}$$

where:

- $P(C_k|X)$ is the posterior probability of class C_k given data X .
- $P(X|C_k)$ is the likelihood of observing data X given class C_k .
- $P(C_k)$ is the prior probability of class C_k .
- $P(X)$ is the probability of data X (a normalization factor).

Multinomial Naive Bayes for Text Classification For text classification, we use the multinomial distribution:

$$P(x_i|C_k) = \frac{N_{x_i,C_k} + \alpha}{N_{C_k} + \alpha d}$$

where:

- N_{x_i,C_k} is the count of word x_i in class C_k .
- N_{C_k} is the total number of words in class C_k
- d is the vocabulary size.
- α is the Laplace smoothing parameter.

6.4.2 Key Parameters and Their Impact

The classifier is tuned using GridSearchCV on the alpha parameter:

```
1 MultinomialNB(alpha=best_alpha)
```

- **alpha:** This is the smoothing parameter (Laplacian smoothing). A small alpha prevents zero probabilities for unseen words, while a large alpha smooths too aggressively.
- The best *alpha* is determined via cross-validation (cv=5), optimizing accuracy.
- After tuning, the final model is trained on the entire dataset with the best alpha value.

6.4.3 Feature Importance Analysis

Since Naive Bayes is a probabilistic model, feature importance can be interpreted based on the learned conditional probabilities $P(X|C)$:

- Words that are strongly associated with a particular class will have higher conditional probabilities in that class.
- The impact of each feature (word) depends on its frequency and how distinctively it appears in a single category.
- In text classification, stopwords and common words contribute less, whereas category-specific words have a stronger influence.

6.4.4 Strengths and Limitations

- **Strengths**

- *Fast and efficient:* Training and inference are computationally inexpensive.
- *Handles high-dimensional data well:* Works effectively for text classification.
- *Interpretable:* Probabilities provide insights into the importance of words.

- **Limitations**

- *Strong independence assumption:* Assumes words appear independently which is often unrealistic in natural language.
- *Performance on complex patterns:* Struggles with tasks where feature dependencies are important.
- *Sensitivity to rare words:* If a word appears very few times, it may not contribute effectively to classification.

6.5 Graphical Model

6.5.1 Overview

Graphical models are probabilistic models that use graphs to represent dependencies between variables, making them well-suited for structured prediction tasks. We conduct our topic on two widely used graphical models: Bayesian Network (BN) and Hidden Markov Model (HMM). These models provide a structured way to model relationships between words and categories, leveraging the dependencies rather than assuming complete independence as in Naive Bayes.

- **Bayesian Network (BN):** A directed acyclic graph (DAG) where nodes represent variables (words, document categories,...) and edges encode probabilistic dependencies. The joint probability distribution over all variables is factorized as:

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | Pa(X_i))$$

where $Pa(X_i)$ represents the parent nodes of X_i in the DAG.

- **Hidden Markov Model (HMM):** A sequential probabilistic model where observations (words) are generated by the hidden states (topics or categories). The probability of an observation sequence $O = (o_1, o_2, \dots, o_T)$ given a hidden state sequence $S = (s_1, s_2, \dots, s_T)$ is computed as:

$$P(O, S) = P(s_1) \prod_{t=2}^T P(s_t | s_{t-1}) \prod_{t=1}^T P(o_t | s_t)$$

where:

- $P(s_1)$ is the initial state probability,
- $P(s_t | s_{t-1})$ is the state transition probability,
- $P(o_t | s_t)$ is the observation likelihood.

6.5.2 Key Parameters and Their Impact

6.5.2.a Bayesian Network

```
1 structure = [('Category', f'word_{i}') for i in range(X_disc.shape[1])]
2 self.model = BayesianNetwork(structure)
3 self.model.fit(df, estimator=MaximumLikelihoodEstimator)
```

1. Structure Learning (BayesianNetwork(structure))

- **Definition:** Structure learning defines the graphical representation of dependencies between variables (features). The structure can be predefined based on domain knowledge or learned from data.
- **Impact on Classification:**
 - A well-designed structure improves classification accuracy by capturing essential feature dependencies.
 - If the structure is too simplistic, the model may miss critical relationships, leading to poor generalization.
 - An overly complex structure may lead to overfitting, making the model less robust to new data.

2. Conditional Probability Tables (CPTs)(fit(df, estimator=MaximumLikelihoodEstimator))

- **Definition:** CPTs store the conditional probabilities of each node given its parent nodes which are estimated from the training data.
- **Estimator choices:**
 - We choose Maximum Likelihood Estimator (MLE) to compute probabilities directly from observed data and it works well when the dataset is large.

3. Inference Algorithm (VariableElimination(self.model))

- **Definition:** Bayesian Networks support different inference techniques, such as exact and approximate methods.
- **Impact:**
 - Variable Elimination: computes exact probabilities by marginalizing out irrelevant variables. Moreover, it is suitable for small to medium-sized models.
 - We use Variable Elimination for inference due to its accuracy. If scalability becomes an issue, approximate methods should be considered.

6.5.2.b Hidden Markov Model

```
1 model = hmm.GaussianHMM(n_components=4, covariance_type="full", tol=1e-3, init_params="
    stmc", verbose=True)
2 model.fit(X_train_scaled[idx])
```

1. Number of Hidden States (n_components)

- **Definition:** Specifies the number of hidden states in the HMM.
- **Impact:** A higher number of components allows the model to capture more complex structures but may lead to overfitting if not properly tuned.

2. covariance_type

- **Definition:** Determines the structure of covariance matrices used in the Gaussian emissions.
- **Options:**
 - "full": Provides the most flexibility but requires more parameters to estimate.
 - "tied": All components share the same covariance matrix.
 - "diag": Each component has a diagonal covariance matrix.
 - "spherical": Each component has a single variance value.
- **Impact:**
 - "full": Provides the most flexibility but requires more parameters to estimate.
 - "diag" and "spherical": Reduce computational complexity at the cost of modeling accuracy.

3. tol

- **Definition:** The convergence threshold for the log-likelihood.
- **Impact:** A smaller `tol` value increases the number of iterations required for convergence, while a higher `tol` may result in premature convergence and suboptimal solutions.

4. verbose

- **Definition:** Controls the verbosity of output during model training.
- **Impact:**
 - `True`: Displays detailed logs, useful for debugging and monitoring convergence.
 - `False`: Suppresses logs, leading to cleaner output during execution.

6.5.3 Feature Importance Analysis

6.5.3.a Bayesian Network

In Bayesian Network, feature importance is inferred from the structure of the network, which represents probabilistic dependencies between variables. In our model, the Category variable is the parent node, and each word feature (word 0, word 1,..., word n) is directly dependent on it. This structure reflects the direct influence of the categorical label on the distribution of words within the dataset. We can evaluate the feature importance based on:

1. **Mutual Information:** Higher mutual information between a word feature and the category suggests that the word is highly indicative of the class.
2. **Conditional Probability Distributions:** The probability of a specific word feature given a category helps identify key discriminative words.
3. **Sensitivity Analysis:** Observing changes in classification outcomes when certain features are perturbed or removed.

Our analysis showed that domain-specific words had higher importance scores, meaning they significantly influenced the classification outcome. For example, words like "*government*" and "*election*" were strongly associated with the *World* category, while "*match*" and "*goal*" were more relevant for *Sports*.

6.5.3.b Hidden Markov Model

HMMs rely on sequential word dependencies. Unlike models use bag-of-words (BoW) or TF-IDF, HMMs assume an underlying state transition structure. The importance of features in HMMs arises from:

- **Word Sequences:** The probability of transitioning from one word to another.
- **Emission Probabilities:** High-probability word sequences help determine class labels.
- **State Transitions:** The likelihood of moving between hidden states influences classification.

If feature selection removes key sequential dependencies, the model loses predictive power.

6.5.4 Strengths and Limitations

6.5.4.a Bayesian Network

- **Strengths**
 - **Interpretable Structure:** The probabilistic dependencies in Bayesian Networks make it easy to understand relationships between features and target categories.
 - **Handle Uncertainty Well:** Unlike deterministic models, Bayesian Networks provide a probabilistic output, allowing confidence estimation for each prediction.
 - **Data Efficiency:** Performs well even with limited training data by leveraging prior probabilities and conditional dependencies.
 - **Resilience to missing data:** Since it models conditional probabilities, it can make reasonable inferences even with partial feature information.
- **Limitations**
 - **Computational Complexity:** Exact inference in Bayesian Networks can be expensive, especially as the number of features increases.
 - **Discretization Sensitivity:** The discretization process in our model (binning continuous TF-IDF features) can lead to information loss if not tuned properly.

6.5.4.b Hidden Markov Model

- **Strengths**
 - Captures sequential relationships in text.
 - Works well for structured language models (POS tagging, speech recognition). Can model probabilistic transitions, unlike purely count-based models.
- **Limitations**
 - **Poor performance in text classification:** Assumes strict sequence dependencies, making it unsuitable for datasets where order does not play a strong role.
 - **High computational cost:** Large vocabularies require intensive matrix operations.
 - **Sensitivity to parameter initialization:** Poor initializations can lead to suboptimal local minima.
 - **Limited generalizability:** Works best in specific domains with strong sequential structures.

7 Usage Example

The trained pipeline is applied to classify new, unseen text samples. This example:

- Demonstrates the practical application of the model in news classification.
- Provides a reference for future use and model extensions.

Example Predictions:

```
new_samples = [  
    "Wall Street sees renewed optimism in tech stocks",  
    "Soccer World Cup final brings excitement worldwide"  
]  
Predicted Label: [  
    "Business",  
    "Sports"  
]
```

This showcases the model's ability to generalize and predict meaningful labels based on text input.

8 Experimental Results

To assess the effectiveness of each model, we will conduct comparative evaluations between models using the same use-case on the introduced datasets across various settings.

8.1 Evaluation Method

For the specific problem of classification based on text, several metrics are used to evaluate the model's performance. In this thesis report, the presenter uses 4 metrics to evaluate the experimental results on our target dataset: *Precision*, *Recall*, *F1 Score*, and *Accuracy*.

The confusion matrix (or error matrix) is a table designed to visualize the results of an algorithm. Because accuracy alone cannot determine the actual performance of classification system if the dataset is unbalanced. To test the detection function and the conditions of the confusion matrix are shown in Table 2.

Condition	Interpretation
True Positive (TP)	Abnormal points correctly detected as actual abnormal points.
True Negative (TN)	Non-abnormal points correctly detected.
False Positive (FP)	Non-abnormal points incorrectly predicted as abnormal.
False Negative (FN)	Abnormal points not detected.

Table 2: Conditions and explanations of values in the confusion matrix

- **Accuracy:** The most intuitive measure, it is simply the ratio of the correctly predicted observations to the total observations. Although widely used, accuracy is not the most appropriate measure in some cases, especially when the classes in the dataset are imbalanced.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

- **F1 Macro score**

The F1 score considers both precision and recall to compute the model's performance. Mathematically, it is the harmonic mean of the model's precision and recall, computed as follows:

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

$$F1 - Score = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (4)$$

This F1 score is computed across different classes, taking the average of the F1 scores of all classes to obtain the average F1 score.

8.2 Experimental Results and Discussion

Rank	Model	Accuracy	Precision	Recall	F1 Score
1	Neural Network	0.8762	0.88	0.88	0.88
2	Naive Bayes	0.8523	0.85	0.85	0.85
3	Bayes Network	0.8433	0.84	0.84	0.84
4	Decision Tree	0.7996	0.80	0.80	0.80
5	HMM	0.6514	0.67	0.65	0.66

Table 3: Results on the **AG News** dataset

Table 3 presents the experimental evaluation results over implemented models and the baseline models on the **AG News** dataset.

8.2.1 Neural Network (MLP)

The Neural Network demonstrated the best classification results across all metrics ($Accuracy = 0.8762$, $Precision = 0.88$, $Recall = 0.88$, $F1\ Score = 0.88$). The MLP classifier benefits from non-linear learning capabilities, which can lead to improved performance on complex tasks such as multi-class text classification (e.g., AG News). Its ability to capture complex, non-linear patterns in text data—especially when coupled with TF-IDF features—contributed to robust generalization. However, despite its strong performance, interpretability remains a challenge. The model's performance is highly dependent on hyperparameter tuning (e.g., learning rate, hidden layer sizes, and regularization) and a fixed number of iterations might not suffice for full convergence on larger datasets.

8.2.2 Naive Bayes

Naive Bayes ($Accuracy = 0.8523$, $Precision = 0.85$, $Recall = 0.85$, $F1\ Score = 0.85$) ranked second, confirming its reputation as a fast and effective baseline for text classification. Naive Bayes is computationally efficient and well-suited for high-dimensional text data. Its probabilistic nature offers a degree of interpretability, as the importance of each word is derived from its learned conditional probability.

Although it assumes feature independence—which is rarely true in natural language—its simplicity and low computational overhead make it appealing for large-scale or real-time classification tasks. However, the strong independence assumption may limit its effectiveness on tasks where word dependencies are critical, and rare words might not contribute effectively. Moreover, the probabilistic framework allows a degree of interpretability by examining the conditional probabilities for each word.

8.2.3 Bayesian Network

The Bayesian Network approach achieved an *Accuracy* of 0.8433 with balanced *Precision*, *Recall*, and *F1 Score* of 0.84. Bayesian Networks excel in interpretability due to their explicit representation of probabilistic dependencies among features and labels. Its key advantage lies in modeling dependencies among words and the target category, providing clearer interpretability. However, learning and inference can become computationally expensive for high-dimensional text data, especially if the network structure grows complex. Discretizing continuous features (e.g., TF-IDF scores) may introduce information loss if not done carefully.

8.2.4 Decision Tree

The Decision Tree classifier attained an *Accuracy* of 0.7996, *Precision* = 0.80, *Recall* = 0.80, and *F1 Score* = 0.80. Decision Trees excel in interpretability, offering transparent, rule-based classifications. While Decision Trees are non-parametric and do not assume any underlying data distribution, they are prone to overfitting, especially when the tree grows too deep. Furthermore, the high dimensionality introduced by the vectorization of TF-IDF can complicate the decision boundaries and be more susceptible to noise and borderline misclassifications, limiting the generalization capacity of the model compared to ensemble approaches such as Random Forest.

8.2.5 Graphical Models

Finally, the HMM yielded the lowest performance (*Accuracy* = 0.6514, *F1 Score* = 0.66). While HMMs are powerful for tasks involving strong sequential dependencies, they are less effective for general news classification, where strict word-order patterns are less critical. Their reliance on state transitions and emission probabilities can lead to higher complexity and reduced scalability in broad-domain text classification. They also require careful parameter initialization and can be computationally expensive for large vocabularies.

8.2.6 Overall Discussion

The results show that no single model is universally superior in every aspect:

- **Neural Network** outperforms others in terms of raw predictive metrics but requires careful tuning and lacks direct interpretability.
- **Naive Bayes** and **Bayesian Network** balance performance, computational cost, and interpretability, making them attractive for many text classification scenarios.
- **Decision Tree** remains valuable for its human-readable rules but is prone to overfitting, especially with sparse, high-dimensional data.
- **Bayesian Network** and **HMM** is less suitable for non-sequential tasks due to its reliance on strict word-order patterns.

9 Conclusion

This project implementation for a 4-label text classification task exemplifies robust machine learning engineering practices. The comprehensive approach—including data preprocessing, model tuning, performance evaluation, and feature importance analysis—aligns well with the project requirements, demonstrating the integration of theoretical knowledge and practical skills.

Furthermore, this project applies the models learned in the Machine Learning course to a specific classification task, ensuring consistency in the use case across different algorithms. By evaluating various models under the same conditions, we assess their effectiveness in handling text classification problems. This structured approach highlights the strengths and limitations of each model, reinforcing key concepts such as interpretability, generalization, and optimization.

Through this implementation, we bridge theoretical learning with real-world application, strengthening our understanding of machine learning techniques and their impact on practical tasks.

10 References

- [1] Hirsh, H. Generalizing version spaces. Machine Learning, 1994
- [2] L.Breiman, J.Friedman, R.Olshen and C.Stone. Classification and Regression Trees, 1984.
- [3] Haykin, S. Neural Networks and Learning Machines 3rd Edition, 2008.
- [4] Devroye, L., Györfi, L., Lugosi. A Probabilistic Theory of Pattern Recognition, 1996.
- [5] Koller, Friedman. Probabilistic Graphical Models: Principles and Techniques, 2009.
- [6] Bishop, C.M. Pattern Recognition and Machine Learning, 2006.
- [7] Hastie, Tibshirani, Friedman. The Elements of Statistical Learning 2nd Edition, 2009.