

AUTHORSHIP ATTRIBUTION, EMOTION AND SENTIMENT PREDICTION IN TAMIL POEMS USING AUTHORSHIP IDENTIFICATION FEATURES

By Team Chozha Pudhalvargal

Internship Report

Submitted by:

Sakthi Aakarshan R
[127145005]

Kishore M
[127177013]

Vennila K M
[127177033]

Harish K
[127177008]

Under the Guidance of:

Prof. Dr.Santhi.B
[Dean SRC Kumbakonam]

SASTRA Deemed University
SASHE

July 6, 2025

Contents

1	Abstract	3
2	Introduction	3
2.1	Background	3
2.2	Problem Statement	3
2.3	Objectives	4
3	Literature Review	4
3.1	Authorship Attribution	4
3.2	Emotion and Sentiment Analysis	4
3.3	Tamil NLP Research	5
4	Methodology	5
4.1	Data Collection	5
4.1.1	Poet Selection	5
4.1.2	Data Structure	5
4.2	Preprocessing and Feature Engineering	6
4.2.1	Preprocessing Steps	6
4.2.2	Feature Engineering	6
4.3	Data Preparation	7
4.3.1	Data Consolidation	7
4.3.2	Manual Annotation	8
4.4	Model Selection and Comparison	8
4.4.1	Algorithm Selection	8
4.4.2	Evaluation Framework	9
4.5	Model Training and Hyperparameter Tuning	9
4.5.1	Training Process	9
4.5.2	Hyperparameter Optimization	9
5	Implementation	10
5.1	Technical Architecture	10
5.1.1	Core Components	10
5.1.2	Technology Stack	10
5.2	Web Application Development	10
5.2.1	Application Features	10
5.2.2	User Interface Design	11

6	Results and Discussion	11
6.1	Model Performance	11
6.1.1	Authorship Attribution Results	11
6.1.2	Emotion Classification Results	11
6.1.3	Sentiment Analysis Results	12
6.2	Feature Importance Analysis	12
6.2.1	Authorship Attribution	12
6.2.2	Emotion Classification	12
6.2.3	Sentiment Analysis	13
6.3	Web Application Performance	13
7	Challenges and Limitations	13
7.1	Technical Challenges	13
7.1.1	Tamil Language Processing	13
7.1.2	Feature Engineering	13
7.2	Dataset Limitations	14
7.2.1	Size Constraints	14
7.2.2	Annotation Challenges	14
7.3	Methodological Limitations	14
7.3.1	Feature Representation	14
7.3.2	Evaluation Scope	14
8	Future Work	15
8.1	Immediate Enhancements	15
8.1.1	Pipeline Integration	15
8.1.2	Model Improvements	15
8.2	Long-term Research Directions	15
8.2.1	Dataset Expansion	15
8.2.2	Advanced Feature Engineering	16
8.2.3	Interdisciplinary Applications	16
9	Conclusion	16
9.1	Key Contributions	16
9.2	Performance Achievements	17
9.3	Significance and Impact	17
9.4	Lessons Learned	18
9.5	Final Remarks	18

1 Abstract

This report presents a comprehensive study on computational analysis of Tamil poetry, focusing on three primary objectives: authorship attribution, emotion classification, and sentiment analysis. The project employs machine learning techniques combined with traditional natural language processing methods to analyze poems from four renowned Tamil poets: Avvaiyar, Kambar, Manikavasagar, and Kapilar. A total of 160 poems were collected and analyzed using 24 distinct authorship identification features, ranging from lexical and syntactic characteristics to poetic structural elements. The study compares ten different machine learning algorithms across all three classification tasks, with Random Forest achieving optimal performance for authorship attribution and sentiment analysis, while Decision Tree proved most effective for emotion classification. The research culminates in a web-based application that provides real-time predictions and visualizations, making the analytical framework accessible for practical applications in Tamil literary studies.

2 Introduction

2.1 Background

Tamil literature represents one of the oldest and richest literary traditions in the world, spanning over two millennia. The computational analysis of Tamil poetry presents unique challenges due to the language's complex morphological structure, rich poetic conventions, and diverse stylistic variations across different periods and authors. Traditional literary analysis, while valuable, is inherently subjective and time-consuming when applied to large corpora.

The advent of computational linguistics and machine learning has opened new avenues for objective literary analysis. Authorship attribution, the process of identifying the author of a text based on stylistic features, has found successful applications in various languages and literary forms. Similarly, emotion and sentiment analysis have become crucial tools for understanding the emotional landscape of literary works.

2.2 Problem Statement

The primary challenge addressed in this project is the development of an automated system capable of:

1. Accurately identifying the authorship of Tamil poems based on stylistic features
2. Classifying emotional content within Tamil poetry

3. Determining sentiment polarity of poetic expressions
4. Providing an accessible interface for practical application of these analyses

2.3 Objectives

The main objectives of this research are:

1. To collect and curate a comprehensive dataset of Tamil poems from classical poets
2. To develop a robust feature extraction pipeline for Tamil poetry analysis
3. To implement and compare multiple machine learning algorithms for multi-task classification
4. To create a web-based application for real-time poetry analysis
5. To establish a foundation for future research in computational Tamil literature analysis

3 Literature Review

3.1 Authorship Attribution

Authorship attribution has been a subject of computational linguistics research for several decades. Early work by Mosteller and Wallace (1964) on the Federalist Papers established the foundation for statistical approaches to authorship identification. Modern techniques have evolved to incorporate sophisticated feature engineering and machine learning algorithms.

Key features commonly used in authorship attribution include:

- Lexical features: word frequencies, vocabulary richness, word length statistics
- Syntactic features: part-of-speech distributions, sentence structures
- Stylistic features: punctuation patterns, function word usage
- Character-level features: n-gram distributions, character entropy

3.2 Emotion and Sentiment Analysis

Emotion analysis in literature has gained significant attention with the development of computational approaches. Russell's circumplex model and Ekman's basic emotions theory provide theoretical foundations for emotion classification systems. In the context of

poetry, emotional analysis becomes more complex due to metaphorical language, cultural references, and artistic expression.

Sentiment analysis, while conceptually simpler than emotion analysis, presents unique challenges in poetic texts where conventional sentiment indicators may not apply directly.

3.3 Tamil NLP Research

Research in Tamil natural language processing has been growing, with contributions in areas such as:

- Morphological analysis and part-of-speech tagging
- Machine translation systems
- Text classification and information retrieval
- Computational poetry analysis (limited but emerging)

4 Methodology

4.1 Data Collection

The dataset was Manually collected from various sources to ensure comprehensive representation of the selected poets' works. The collection process involved:

4.1.1 Poet Selection

Four renowned Tamil poets were selected based on their historical significance and distinct poetic styles:

1. **Avvaiyar**: Known for moral and philosophical poetry
2. **Kambar**: Famous for epic poetry, particularly the Ramayana
3. **Manikavasagar**: Devotional poetry and spiritual themes
4. **Kapilar**: Classical Sangam period poetry

4.1.2 Data Structure

The dataset was organized hierarchically:

- Total poems: 160 (40 poems per poet)
- Directory structure: 4 folders named after each poet

- File format: Individual text files (.txt) for each poem
- Naming convention: Poem titles as filenames

4.2 Preprocessing and Feature Engineering

A comprehensive Python pipeline was developed to handle preprocessing and feature extraction tasks. The pipeline architecture ensures consistency and reproducibility across all data processing stages.

4.2.1 Preprocessing Steps

1. **Text Loading:** Systematic reading of all poem files from the organized directory structure
2. **Tokenization:** Breaking down poems into individual words and characters for analysis
3. **Normalization:** Handling variations in text encoding and formatting
4. **Feature Extraction:** Computing 24 distinct authorship identification features

4.2.2 Feature Engineering

The feature set was designed to capture multiple aspects of poetic style and structure:

Lexical Features:

- TOTAL_WORD_COUNT: Total number of words in the poem
- UNIQUE_WORD_COUNT: Number of unique words used
- TYPE_TOKEN_RATIO: Ratio of unique words to total words
- AVG_WORD_LENGTH: Average length of words in characters
- STD_WORD_LENGTH: Standard deviation of word lengths
- HAPAX_LEGOMENA: Words that appear only once

Syntactic Features:

- FUNCTION_WORD_FREQ: Frequency of function words
- STOP_WORD_RATIO: Ratio of stop words to total words
- ESTIMATED_NOUN_COUNT: Estimated number of nouns
- ESTIMATED_VERB_COUNT: Estimated number of verbs

- NOUN_VERB_RATIO: Ratio of nouns to verbs

Character-level Features:

- CHAR_UNIGRAM_ENTROPY: Entropy of character unigrams
- CHAR_BIGRAM_ENTROPY: Entropy of character bigrams
- CHAR_TRIGRAM_ENTROPY: Entropy of character trigrams
- CHAR_DIVERSITY: Diversity of character usage
- RARE_CHAR_RATIO: Ratio of rarely used characters

Poetic Structure Features:

- AVG_SYLLABLES_PER_WORD: Average syllables per word
- AVG_LINE_LENGTH_WORDS: Average words per line
- AVG_STANZA_LENGTH: Average lines per stanza
- TOTAL_LINES: Total number of lines
- ANAPHORA_SCORE: Repetition at line beginnings
- EPIPHORA_SCORE: Repetition at line endings
- ENJAMBMENT_RATIO: Ratio of enjambed lines

Readability Features:

- READABILITY_SCORE: Computational readability measure

4.3 Data Preparation

4.3.1 Data Consolidation

Individual CSV files generated for each poet were concatenated to create a unified dataset. This process involved:

1. Merging CSV files while preserving author labels
2. Ensuring consistent feature representation across all entries
3. Validating data integrity and completeness

4.3.2 Manual Annotation

Emotion and sentiment labels were manually assigned to each poem through collaborative analysis by the research team. This process involved:

1. Careful reading and interpretation of each poem
2. For emotion classification, three classes (happy, sad, devotion) were used.
3. For sentiment classification, 3 classes (positive, neutral, Negative) were used
4. Quality assurance through multiple reviewer validation

The manual annotation process ensures high-quality ground truth labels, though it represents a time-intensive aspect of the methodology.

4.4 Model Selection and Comparison

A comprehensive evaluation of machine learning algorithms was conducted to identify optimal models for each classification task.

4.4.1 Algorithm Selection

Ten traditional machine learning algorithms were selected for comparison:

1. Logistic Regression
2. Random Forest
3. Support Vector Machine (SVM)
4. K-Nearest Neighbors (KNN)
5. Naive Bayes
6. Decision Tree
7. Neural Network
8. XGBoost
9. LightGBM
10. Gradient Boosting

4.4.2 Evaluation Framework

Each algorithm was evaluated across three distinct classification tasks:

1. **Authorship Attribution:** 4-class classification (one per poet)
2. **Emotion Classification:** Multi-class emotion categorization
3. **Sentiment Analysis:** Multi-class sentiment classification

4.5 Model Training and Hyperparameter Tuning

4.5.1 Training Process

The best-performing models identified during the comparison phase were subjected to comprehensive hyperparameter optimization:

Authorship Attribution - Random Forest:

- Optimized parameters: `n_estimators`, `max_depth`, `min_samples_split`
- Performance metric: Accuracy

Emotion Classification - Decision Tree:

- Optimized parameters: `criterion`, `max_depth`, `min_samples_leaf`
- Performance metric: Accuracy

Sentiment Analysis - Random Forest:

- Optimized parameters: `n_estimators`, `max_features`, `bootstrap`
- Performance metric: Accuracy

4.5.2 Hyperparameter Optimization

Grid search methodology was employed to systematically explore hyperparameter spaces:

- `test_size`: [0.2, 0.23, 0.25, 0.28, 0.3, 0.33, 0.35, 0.38, 0.4]
- `dataset_random_state`: 1 to 100
- `model_random_state`: [42, 44, 46, 47]
- `n_estimators`: 100 to 5000
- `criterion`: {'gini', 'entropy', 'log_loss'}
- `max_depth`: 3 to 30

5 Implementation

5.1 Technical Architecture

The project implementation follows a modular architecture designed for scalability and maintainability:

5.1.1 Core Components

1. **Data Processing Pipeline:** Python-based preprocessing and feature extraction
2. **Machine Learning Module:** Model training, evaluation, and prediction
3. **Web Application:** Flask-based web interface
4. **Visualization Module:** Word cloud generation

5.1.2 Technology Stack

- **Backend:** Python, Flask framework
- **Frontend:** HTML5, CSS3, JavaScript
- **Machine Learning:** scikit-learn, pandas, numpy
- **Data Processing:** pandas

5.2 Web Application Development

A user-friendly web interface was developed to make the analytical capabilities accessible to researchers and enthusiasts:

5.2.1 Application Features

1. **Poem Input:** Text area for entering Tamil poems
2. **Multi-task Prediction:** Radio buttons for each task (authorship, emotion, and sentiment prediction)
3. **Word Cloud Visualization:** Dynamic word cloud generation
4. **Results Display:** Comprehensive prediction results

5.2.2 User Interface Design

The web interface was designed with the following principles:

- Responsive design for various screen sizes
- Intuitive navigation and user flow
- Clear presentation of results
- Accessibility considerations for Tamil text display

6 Results and Discussion

6.1 Model Performance

The comparative analysis of machine learning algorithms revealed distinct patterns of performance across the three classification tasks:

6.1.1 Authorship Attribution Results

Random Forest emerged as the optimal algorithm for authorship attribution:

- **Accuracy:** 98.21%
- **Precision:** 0.98
- **Recall:** 0.98
- **F1-Score:** 0.98

The high performance can be attributed to the ensemble nature of Random Forest, which effectively captures the complex stylistic variations between poets while maintaining robustness against overfitting.

6.1.2 Emotion Classification Results

Decision Tree demonstrated superior performance for emotion classification:

- **Accuracy:** 86.49%
- **Precision:** 0.88
- **Recall:** 0.86
- **F1-Score:** 0.87

The interpretability of Decision Tree models proves valuable for understanding the decision-making process in emotion classification, allowing for better insights into the relationship between textual features and emotional content.

6.1.3 Sentiment Analysis Results

Random Forest again proved optimal for sentiment analysis:

- **Accuracy:** 87.5%
- **Precision:** 0.58
- **Recall:** 0.58
- **F1-Score:** 0.58

The multi-class sentiment classification, combined with the robust feature representation, contributed to the strong performance of Random Forest in this task.

6.2 Feature Importance Analysis

Analysis of feature importance revealed interesting patterns in the contribution of different feature types:

6.2.1 Authorship Attribution

Key discriminative features for authorship attribution:

1. Character-level entropy measures (CHAR_TRIGRAM_ENTROPY)
2. Lexical diversity measures (TYPE_TOKEN_RATIO)
3. Poetic structure features (AVG_STANZA_LENGTH)
4. Syntactic patterns (NOUN_VERB_RATIO)

6.2.2 Emotion Classification

Important features for emotion classification:

1. Readability measures (READABILITY_SCORE)
2. Lexical complexity (AVG_WORD_LENGTH)
3. Structural features (TOTAL_LINES)
4. Character diversity (CHAR_DIVERSITY)

6.2.3 Sentiment Analysis

Significant features for sentiment analysis:

1. Function word frequency (FUNCTION_WORD_FREQ)
2. Poetic devices (ANAPHORA_SCORE, EPIPHORA_SCORE)
3. Lexical measures (HAPAX_LEGOMENA)
4. Structural complexity (ENJAMBMENT_RATIO)

6.3 Web Application Performance

The web application demonstrates satisfactory performance characteristics:

- **Response Time:** Average 2 seconds for each analysis
- **Accuracy:** Maintains model accuracy in real-time predictions
- **User Experience:** Intuitive interface with clear result presentation
- **Scalability:** Capable of handling multiple concurrent users

7 Challenges and Limitations

7.1 Technical Challenges

7.1.1 Tamil Language Processing

Processing Tamil text presents unique challenges:

- Complex morphological structure requiring specialized tokenization
- Encoding inconsistencies across different sources
- Limited availability of Tamil NLP resources
- Difficulty in accurate part-of-speech tagging

7.1.2 Feature Engineering

Several challenges were encountered in feature extraction:

- Adaptation of existing features to Tamil linguistic characteristics
- Handling of poetic conventions specific to Tamil literature
- Balancing feature comprehensiveness with computational efficiency

7.2 Dataset Limitations

7.2.1 Size Constraints

The current dataset, while comprehensive for the selected poets, has limitations:

- Relatively small sample size (160 poems total)
- Limited representation of Tamil poetry's full diversity
- Potential bias towards classical poetry forms

7.2.2 Annotation Challenges

Manual annotation of emotion and sentiment presents inherent limitations:

- Subjectivity in emotion interpretation
- Cultural context dependencies
- Time-intensive process limiting scalability

7.3 Methodological Limitations

7.3.1 Feature Representation

Current feature engineering approaches have limitations:

- Reliance on hand-crafted features rather than learned representations
- Limited capture of semantic relationships
- Potential loss of contextual information

7.3.2 Evaluation Scope

The evaluation framework, while comprehensive, has constraints:

- Focus on accuracy metrics without detailed error analysis
- Limited cross-validation with other Tamil poetry corpora
- Absence of baseline comparisons with human annotators

8 Future Work

8.1 Immediate Enhancements

8.1.1 Pipeline Integration

The next phase of development focuses on creating a unified pipeline that integrates all components:

- End-to-end automation from data input to result generation
- Streamlined workflow for batch processing
- Enhanced error handling and logging capabilities
- Configuration management for different use cases

8.1.2 Model Improvements

Several enhancements are planned for model performance:

- Implementation of deep learning approaches (LSTM, BERT)
- Exploration of ensemble methods combining multiple algorithms
- Integration of semantic embeddings for Tamil text
- Development of transfer learning approaches

8.2 Long-term Research Directions

8.2.1 Dataset Expansion

Future work will focus on expanding the dataset:

- Inclusion of additional poets from different periods
- Incorporation of contemporary Tamil poetry
- Development of automated annotation techniques
- Creation of benchmark datasets for Tamil poetry analysis

8.2.2 Advanced Feature Engineering

Research into more sophisticated feature representations:

- Graph-based features capturing structural relationships
- Phonetic features specific to Tamil prosody
- Cultural and thematic feature extraction
- Multi-modal features incorporating meter and rhythm

8.2.3 Interdisciplinary Applications

Exploration of broader applications:

- Digital humanities research platforms
- Educational tools for Tamil literature study
- Historical linguistic analysis
- Cultural preservation initiatives

9 Conclusion

This research presents a comprehensive approach to computational analysis of Tamil poetry, successfully addressing three fundamental challenges in literary analysis: authorship attribution, emotion classification, and sentiment analysis. The project demonstrates the feasibility and effectiveness of applying machine learning techniques to Tamil literary texts while respecting the unique characteristics of the language and its poetic traditions.

9.1 Key Contributions

The major contributions of this work include:

1. **Comprehensive Dataset:** Creation of a well-structured dataset of 160 Tamil poems from four classical poets, providing a foundation for future research in computational Tamil literature analysis.
2. **Feature Engineering Framework:** Development of a robust set of 24 authorship identification features specifically adapted for Tamil poetry, capturing lexical, syntactic, structural, and poetic characteristics.

3. **Multi-task Classification System:** Implementation and evaluation of a multi-task classification framework capable of simultaneously addressing authorship attribution, emotion classification, and sentiment analysis.
4. **Comparative Analysis:** Systematic comparison of ten machine learning algorithms across all three tasks, providing insights into the most effective approaches for each classification challenge.
5. **Practical Application:** Development of a web-based application that makes the analytical capabilities accessible to researchers, educators, and enthusiasts of Tamil literature.

9.2 Performance Achievements

The project achieved notable performance results:

- Authorship attribution accuracy of 98.21% using Random Forest
- Emotion classification accuracy of 86.47% using Decision Tree
- Sentiment analysis accuracy of 87.5% using Random Forest

These results demonstrate the effectiveness of the proposed approach and establish a strong baseline for future research in this domain.

9.3 Significance and Impact

The work contributes to multiple areas of research and application:

Academic Research: The project provides a methodological framework and baseline results that can inform future research in computational Tamil literature analysis, digital humanities, and cross-linguistic authorship attribution studies.

Cultural Preservation: By developing computational tools for Tamil poetry analysis, the work contributes to the preservation and promotion of Tamil literary heritage, making it more accessible to contemporary audiences.

Educational Applications: The web-based interface can serve as an educational tool for students and researchers studying Tamil literature, providing objective insights into poetic characteristics and authorial styles.

Technological Innovation: The project demonstrates the successful adaptation of computational linguistic techniques to a morphologically rich language, contributing to the broader field of natural language processing for Indian languages.

9.4 Lessons Learned

Several important insights emerged from this research:

1. **Feature Engineering Importance:** The success of traditional machine learning approaches highlights the continued importance of domain-specific feature engineering, even in the era of deep learning.
2. **Multi-task Considerations:** Different classification tasks benefit from different algorithmic approaches, suggesting the value of task-specific optimization rather than one-size-fits-all solutions.
3. **Manual Annotation Value:** Despite the time-intensive nature of manual annotation, the high-quality ground truth labels it provides remain crucial for supervised learning approaches.
4. **Cultural Context Significance:** The analysis of Tamil poetry requires careful consideration of cultural and linguistic contexts, emphasizing the need for domain expertise in computational literary analysis.

9.5 Final Remarks

This internship project represents a successful integration of traditional literary scholarship with modern computational techniques, demonstrating the potential for technology to enhance rather than replace human understanding of literature. The work establishes a foundation for continued research in computational analysis of Tamil literature while providing practical tools for immediate application.

The journey from data collection to web application deployment has provided valuable experience in the complete lifecycle of a machine learning project, from problem formulation and data preparation to model development and deployment. The challenges encountered and solutions developed contribute to the growing body of knowledge in applying computational methods to literary analysis, particularly for under-resourced languages like Tamil.

As we look toward the future, the integration of advanced deep learning techniques, expanded datasets, and enhanced feature representations promises to further advance the field of computational Tamil literature analysis, opening new possibilities for understanding and preserving one of the world's great literary traditions.

Acknowledgments

We would like to express our sincere gratitude to Prof. Dr.Santhi.B for her invaluable guidance and support throughout this internship project. Special thanks to our team

who contributed to the manual annotation process and provided valuable insight into Tamil literary analysis. We also acknowledge the various sources from which the poetry dataset was compiled and the open-source community whose tools and libraries made this research possible.