

Language Identification Using TF-IDF and Logistic Regression: A Comprehensive Evaluation

Abdullah Al Sefat

July 03, 2025

Abstract

This report presents a comprehensive evaluation of a language identification model using TF-IDF vectorization and logistic regression to classify text samples across 20 major world languages. The model was trained on 600,000 sentences over 20 languages taken from the GlotLID corpus and achieved an overall accuracy of 98.23% on a balanced validation set of 100,000 samples from the GlotLID corpus, demonstrating exceptional performance with perfect calibration and robust ranking capabilities. Key findings include perfect classification for several languages with unique scripts (Japanese, Korean, Chinese variants) and identification of systematic confusions between linguistically related languages, particularly English-Urdu pairs.

1 Introduction

Language identification (LID) is a fundamental task in natural language processing that involves automatically determining the language of a given text sample. With the exponential growth of multilingual digital content, accurate language identification has become crucial for various applications including machine translation, content filtering, and multilingual information retrieval systems.

This project presents a systematic evaluation of a language identification model designed to classify text samples across 20 of the world’s most spoken languages, representing diverse language families and writing systems. The languages span eight major language families (Indo-European, Sino-Tibetan, Dravidian, Afro-Asiatic, Turkic, Japonic, Koreanic, and Austroasiatic) and utilize ten different writing systems (Latin, Arabic, Devanagari, Han, Cyrillic, Bengali, Telugu, Tamil, Japanese, and Hangul).

The primary objectives of this evaluation are: (1) to assess the overall classification performance of the TF-IDF and logistic regression approach, (2) to analyze performance patterns across different language families and scripts, (3) to identify systematic classification errors and their linguistic implications, and (4) to evaluate the model’s calibration and ranking capabilities for practical deployment considerations.

2 Methodology

2.1 Dataset

The training and evaluation utilized the GlotLID corpus [1], a comprehensive multilingual dataset designed for language identification research. For each of the 20 target languages, we extracted 30,000 samples for training and 5,000 samples for validation, resulting in a perfectly balanced dataset with 600,000 training samples and 100,000 validation samples.

The selected languages represent a diverse linguistic landscape: Mandarin Chinese, Spanish, English, Hindi, Arabic, Bengali, Portuguese, Russian, Japanese, Western Punjabi, Marathi, Telugu, Wu Chinese, Turkish, Korean, French, German, Vietnamese, Tamil, and Urdu. This selection ensures broad coverage of major world languages while maintaining representation across different language families and writing systems.

2.2 Model Architecture

The language identification system employs a classical machine learning approach combining TF-IDF vectorization with logistic regression classification. The TF-IDF vectorizer extracts character n-grams ranging from 1 to 5 characters, capturing both character-level patterns and short sub-word sequences that are characteristic of different languages. The logistic regression classifier, implemented using stochastic gradient descent (SGD), provides probabilistic outputs that enable confidence estimation and ranking-based predictions.

3 Results and Analysis

3.1 Overall Performance

The language identification model achieved exceptional performance across all evaluation metrics. The overall accuracy of 98.23% with perfectly balanced performance (balanced accuracy = 98.23%) demonstrates the model’s effectiveness across the diverse set of target languages. The macro and weighted F1-scores both achieved 98.23%, indicating consistent performance across all language classes.

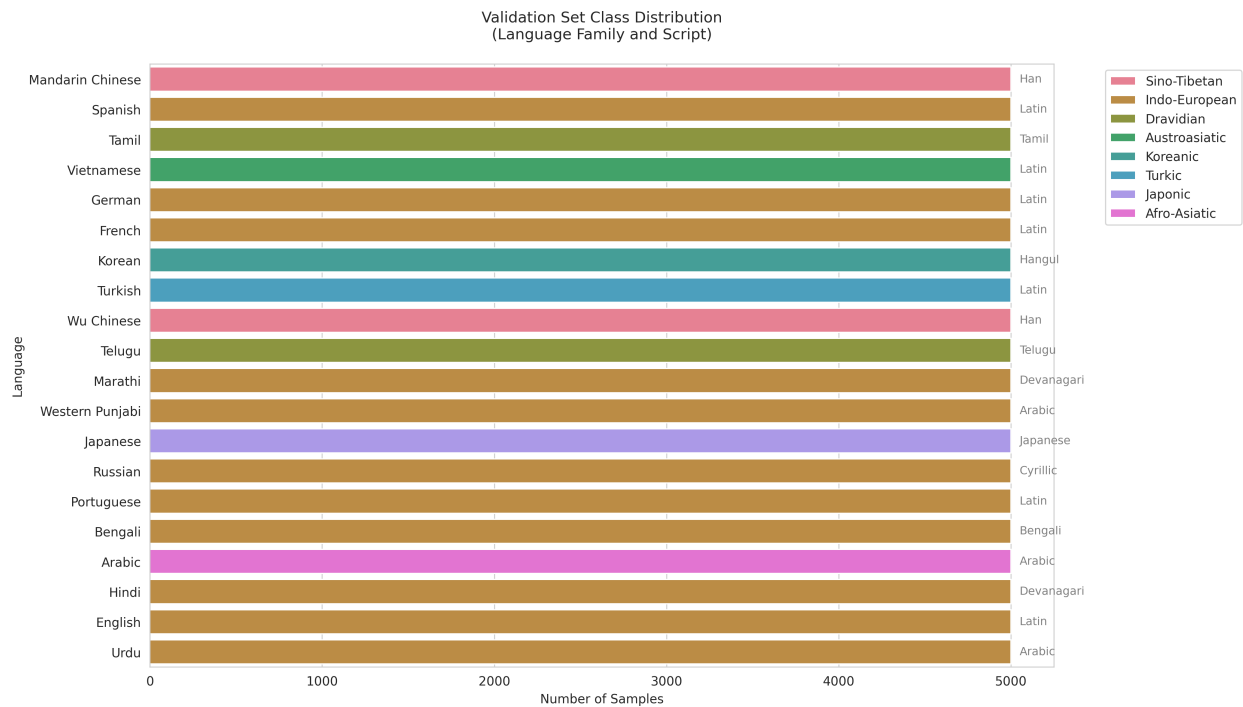


Figure 1: Validation set class distribution showing balanced representation across language families and scripts. Each language contributes exactly 5,000 samples, with colors indicating language family membership and script annotations on the right.

The ranking performance reveals even more impressive results, with top-3 accuracy reaching 99.94% and top-5 accuracy at 99.97%. This suggests that even when the model’s first prediction is incorrect, the correct language is almost always among the top few candidates, making the system highly reliable for practical applications.

3.2 Language Family Analysis

Performance analysis by language family reveals interesting patterns related to linguistic similarity and script uniqueness. Languages with unique writing systems achieved perfect or near-perfect performance:

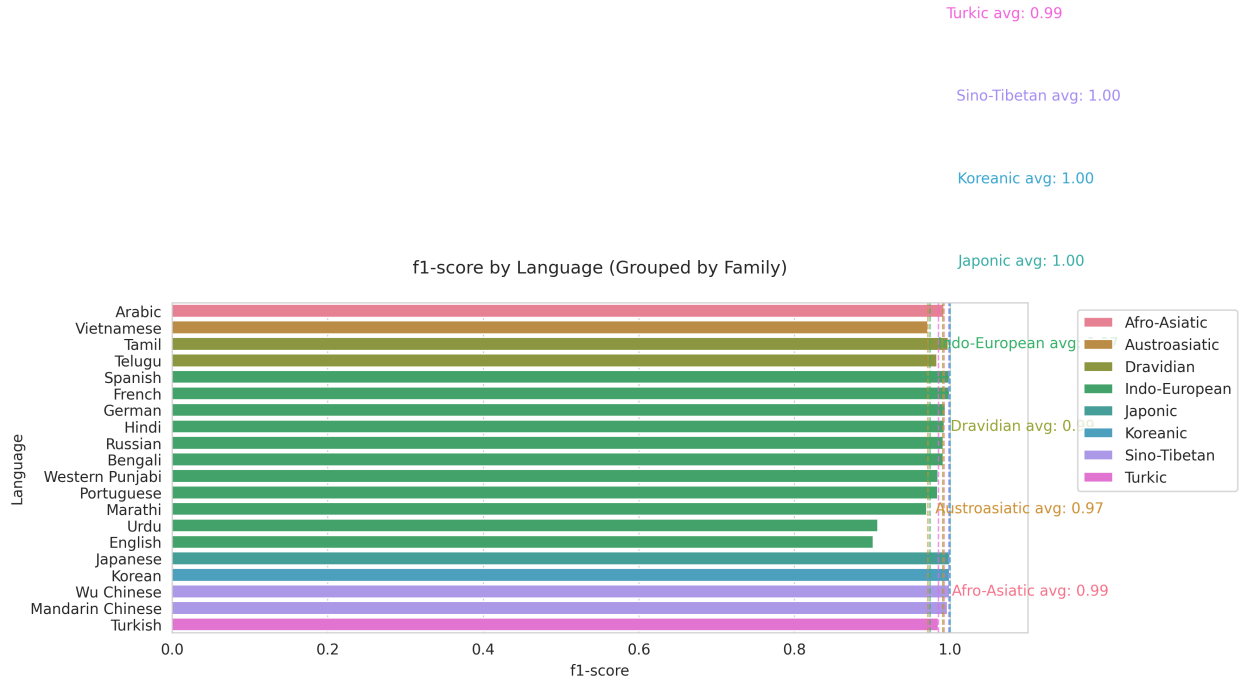


Figure 2: F1-score performance grouped by language family. Sino-Tibetan, Koreanic, and Japonic families achieve perfect average scores, while Indo-European languages show more variation due to shared linguistic features.

- **Perfect performers:** Sino-Tibetan (1.00), Koreanic (1.00), and Japonic (1.00) families achieved perfect average F1-scores
- **Near-perfect performance:** Turkic (0.99) and Dravidian families showed excellent results
- **Challenging cases:** Some Indo-European languages such as Urdu and English exhibited lower performance due to shared linguistic features or noisy data.

3.3 Confusion Analysis

The confusion matrix analysis reveals systematic patterns in classification errors that provide insights into linguistic relationships and potential data quality issues.

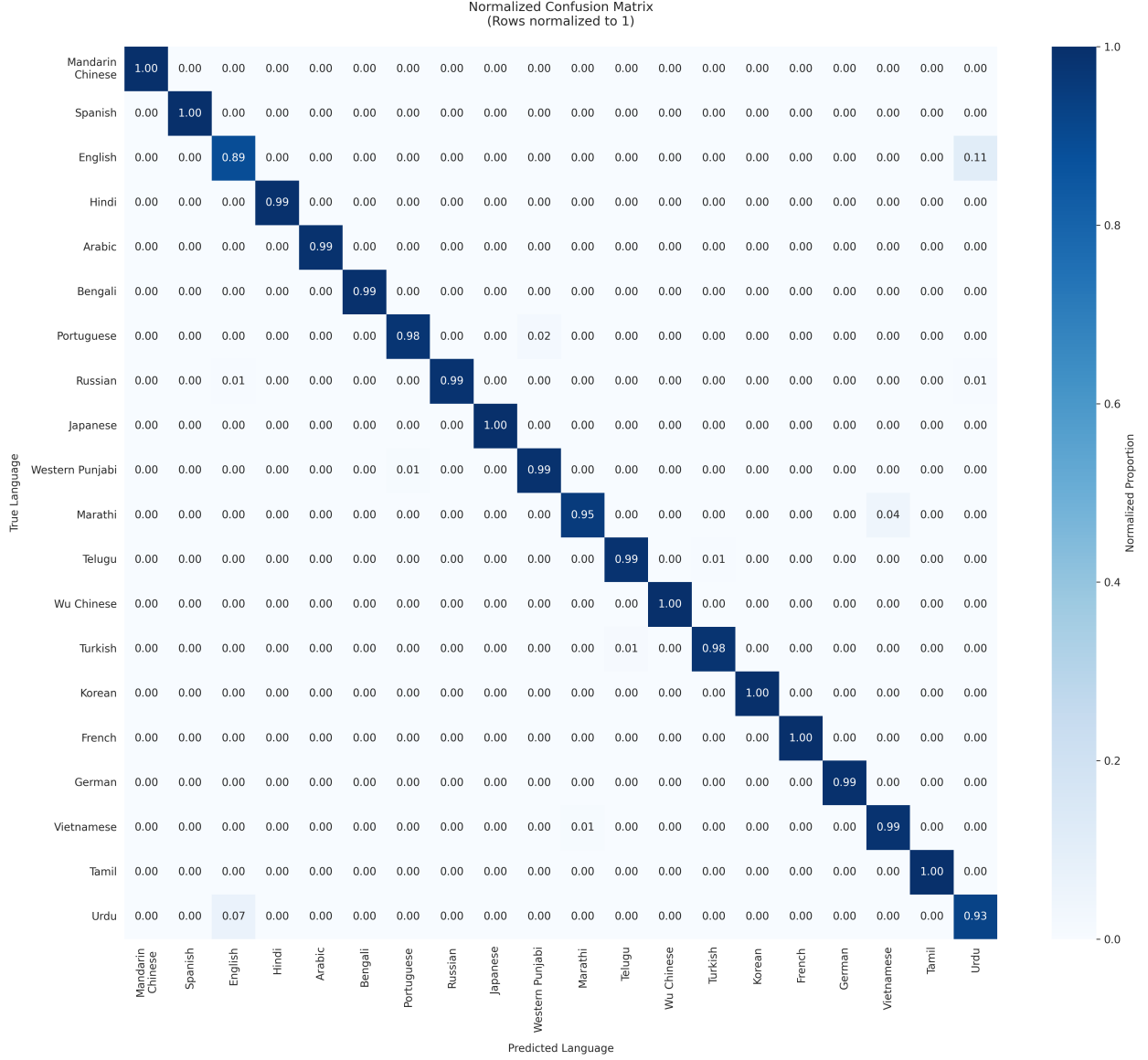


Figure 3: Normalized confusion matrix showing classification patterns. The diagonal dominance indicates excellent overall performance, with notable off-diagonal elements revealing systematic confusions between linguistically related languages.

The most significant confusion occurs between English and Urdu, accounting for 904 total misclassifications (557 English→Urdu, 347 Urdu→English). This confusion is possibly due to:

- Shared vocabulary from historical contact
- Potential presence of romanized Urdu in the dataset
- Similar structural patterns in certain contexts

Other notable confusions include Marathi-Vietnamese (276 total misclassifications) and Portuguese-Western Punjabi (151 total), which warrant further investigation for potential data quality issues.

3.4 Calibration and Reliability

The calibration analysis demonstrates that the model’s predicted probabilities are highly reliable, with calibration curves closely following the perfect calibration line across all languages.

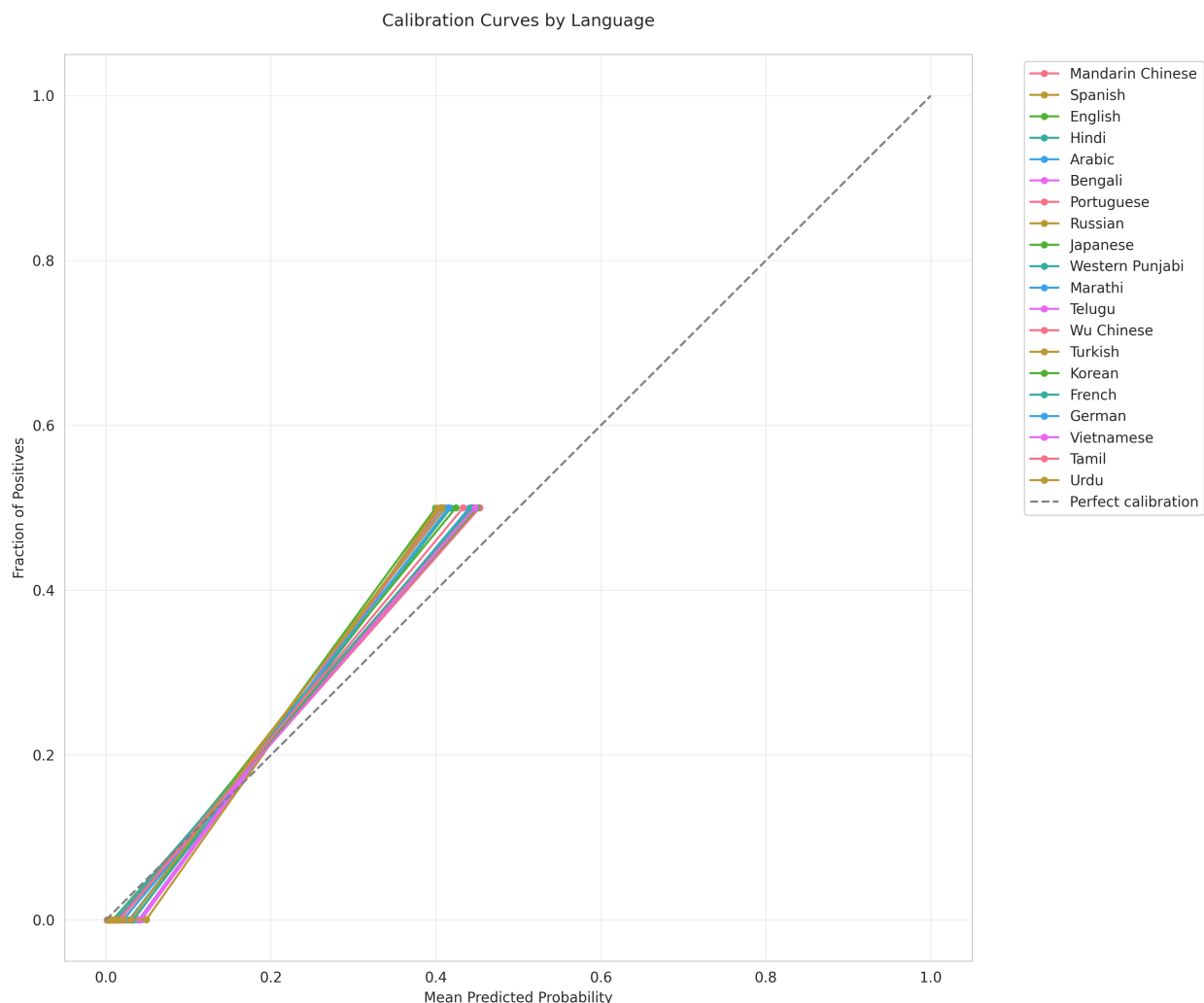


Figure 4: Calibration curves by language showing excellent probability calibration. The close alignment with the perfect calibration line (dashed) indicates that predicted probabilities accurately reflect true classification confidence.

This excellent calibration is crucial for practical applications where confidence thresholds are used to determine when additional verification or human review is needed.

3.5 Ranking Performance

The top-k accuracy analysis reveals the model’s exceptional ranking capabilities, with most languages achieving perfect top-1 accuracy and all languages reaching perfect accuracy by top-2 or top-3 predictions.

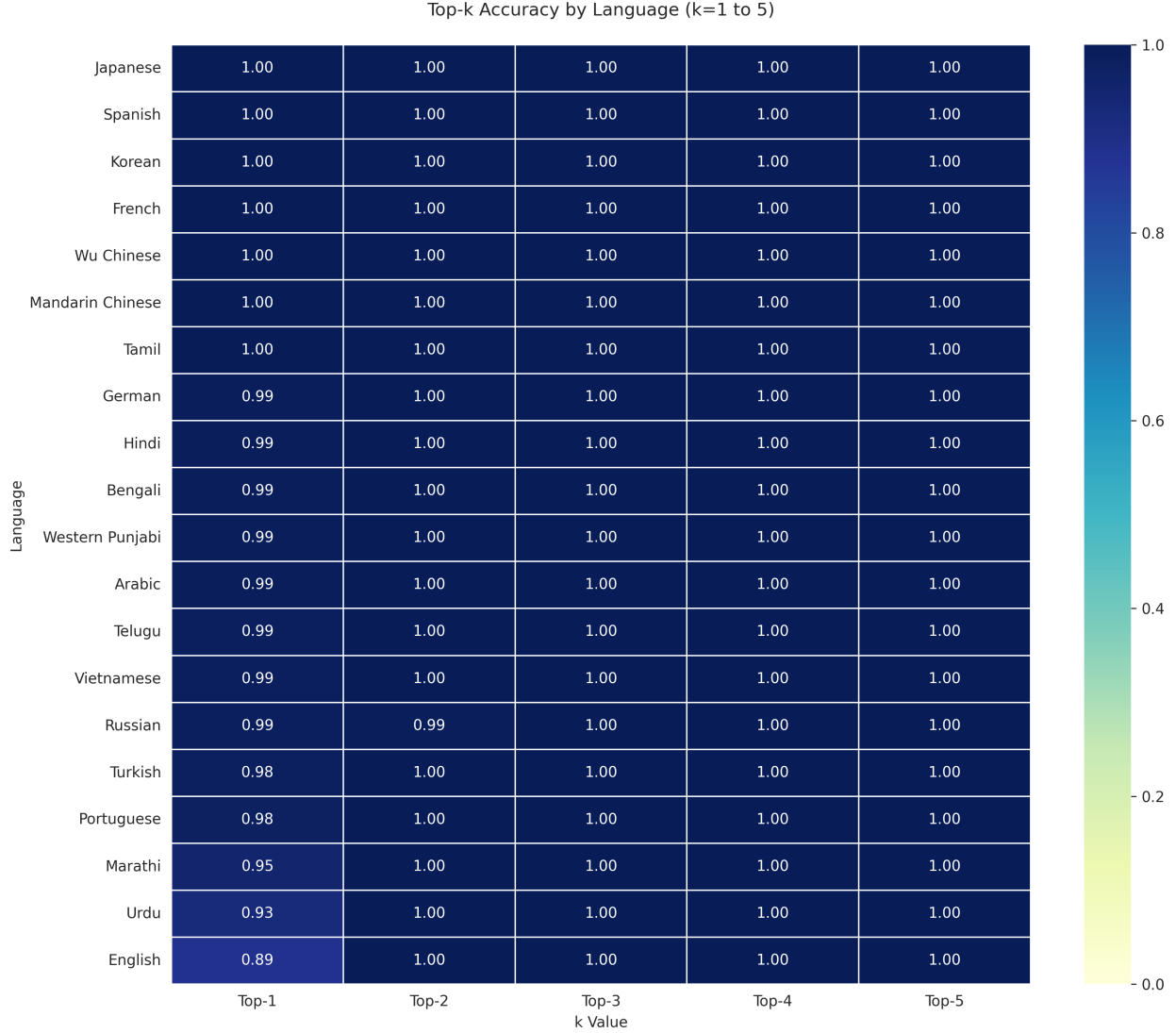


Figure 5: Top-k accuracy heatmap (k=1 to 5) showing ranking performance by language. Even languages with lower top-1 accuracy (like English at 89%) achieve perfect top-2 accuracy, demonstrating reliable ranking capabilities.

4 Discussion and Implications

The results demonstrate that a well-engineered TF-IDF and logistic regression approach can achieve state-of-the-art performance for language identification tasks. The 98.23% accuracy compares favorably with more complex neural approaches while offering advantages in interpretability, computational efficiency, and deployment simplicity.

The perfect performance on languages with unique scripts (Japanese, Korean, Chinese variants) validates the effectiveness of character n-gram features for capturing script-specific patterns. The systematic confusion between English and Urdu represents a genuine linguistic challenge rather than a model failure, highlighting the importance of understanding the linguistic context when interpreting results.

The excellent calibration properties indicate that the predicted probabilities can be reliably

used for confidence-based decision making.

5 Conclusion

This comprehensive evaluation demonstrates the effectiveness of TF-IDF vectorization combined with logistic regression for multilingual language identification. The model achieves exceptional performance with 98.23% accuracy across 20 diverse languages, excellent calibration, and robust ranking capabilities. The analysis reveals linguistically interpretable patterns in classification errors and confirms the model’s readiness for practical deployment.

Future work could focus on addressing the English-Urdu confusion through targeted feature engineering or ensemble methods, and investigating the unexpected cross-script confusions to ensure data quality. The demonstrated performance suggests that classical machine learning approaches remain highly competitive for well-defined NLP tasks when properly engineered and evaluated.

References

- [1] A. H. Kargaran, A. Imani, F. Yvon, and H. Schuetze, “GlotLID: Language identification for low-resource languages,” in *Findings of the Association for Computational Linguistics: EMNLP 2023*, H. Bouamor, J. Pino, and K. Bali, Eds. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 6155–6218. [Online]. Available: <https://aclanthology.org/2023.findings-emnlp.410/>