

EC9640 - ARTIFICIAL INTELLIGENCE
ASSIGNMENT 02

HULANGAMUWA S.C. – 2020/E/054

WIJESINGHE L.A.D.P.A. - 2020/E/182

SEMESTER 07

06 JAN 2024

1. Introduction

This project aims to develop a Sinhala spell and grammar checker that detects and corrects spelling and grammatical errors in Sinhala text. It includes a spell checker for automatic error correction and a grammar checker to identify and fix contextual grammatical issues. Three AI approaches such as Rule-based, LLM-based, and Retrieval-Augmented Generation with Fine-tuned LLM(XLM-RoBERTa) models were evaluated for grammar correction. The most effective approach was selected and tested on five paragraphs to measure its accuracy in spelling corrections and grammar suggestions.

2. Spell correction

a. Description of each considered spelling mistakes their correction approaches with example

Description:

The Sinhala Grammar Checker addresses spelling mistakes like typographical errors, phonetic similarities, and affix variations. Typographical errors are resolved using dictionary matching and edit distance techniques, while phonetic similarities are corrected through advanced phonetic mapping. Prefix and suffix errors are handled with affix rules and stemming to identify root words. The checker uses algorithms like Levenshtein distance, phonetic key generation, and dictionary validation to provide accurate corrections.

Example:

```
--- Spell Check for: භූද්ධාමිත්
Spelling Errors: {'භූද්ධාමිත්': [(('භූද්ධාමිත්', 76.16323529411764), ('භූද්ධාමිත්', 76.16323529411764), ('භූද්ධාමිත්', 75.71323529411765), ('භූද්ධාමිත්', 72.76944444444445), ('භූද්ධාමිත්', 70.41944444444444)]}
Auto-corrected Text: භූද්ධාමිත්

--- Spell Check for: ගාහණ්ථ
Spelling Errors: {'ගාහණ්ථ': [(('ගාහණ්ථ', 82.38333333333333)]}
Auto-corrected Text: ගාහණ්ථ

--- Spell Check for: පනිවිට්ඨ
Spelling Errors: {'පනිවිට්ඨ': [(('පනිවිට්ඨ', 79.90833333333333), ('පිවිට්ඨ', 71.50681818181819), ('නිවිට්ඨ', 71.075)]}
Auto-corrected Text: පනිවිට්ඨ

--- Spell Check for: වානිජ්
Spelling Errors: {'වානිජ්': [(('වානිජ්', 78.025), ('වානිජ්', 71.25)]}
Auto-corrected Text: වානිජ්

--- Spell Check for: මයිකා
Spelling Errors: {'මයිකා': [(('මයිකා', 77.67857142857143), ('මයිකා', 77.67857142857143), ('මයිකා', 70.42500000000001), ('මයිකා', 70.42500000000001)]}
Auto-corrected Text: මයිකා
```

3. Grammar suggestions

a. Description of each considered grammar mistakes (2 different), their suggestions Description:

The Sinhala Grammar Checker code addresses two subject-verb agreement mistakes:

- For "මම" (I), the verb should end with "මි," as in correcting "මම බත් කමු" (I eat rice) to "මම බත් කමි."
- For "අපි" (We), the verb should end with "මු," as in changing "අපි පාඩම් කරමි" (We study) to "අපි පාඩම් කරමු."

The code uses a rule-based approach with predefined suffix mappings, tokenization, and POS tagging to identify and correct these errors, ensuring proper subject-verb agreement.

Example:

```
Processing Input Sentences...

Processing: මම ගදර යයි
Text after Spell Checking: මම ගෙදර යයි
Running Grammar Checker...

Processing: අපි කෑඩට යයි
Text after Spell Checking: අපි කෑඩ යයි
Running Grammar Checker...

Processing: අපි අදරය කරයි
Text after Spell Checking: අපි ආදරය කරයි
Running Grammar Checker...

Processing: මම කෑම ගන්නි
Text after Spell Checking: මම කෑම ගන්නි
Running Grammar Checker...

Processing: අපි රට යයි
Text after Spell Checking: අපි රට යයි
Running Grammar Checker...

Processing: අපි කලී යම
Text after Spell Checking: අපි කලින් යම
Running Grammar Checker...

Final Output:
Sentence 1: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'මි' when subject is 'මම'
Suggested correction: මම ගෙදර යමි
Sentence 2: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'නි' when subject is 'අපි'
Suggested correction: අපි කෑඩ යනි
Sentence 3: Grammar error: Verb 'කරයි' (base: කර, affix: යි) should end with 'නි' when subject is 'අපි'
Suggested correction: අපි ආදරය කරනි
Sentence 4: Grammar error: Verb 'ගන්නි' (base: ගනි, affix: නි) should end with 'මි' when subject is 'මම'
Suggested correction: මම කෑම ගනිමි
Sentence 5: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'නි' when subject is 'අපි'
Suggested correction: අපි රට යනි
Sentence 6: Grammar error: Verb 'යම' (base: ය, affix: ම) should end with 'නි' when subject is 'අපි'
Suggested correction: අපි කලින් යනි
```

b. Existing 3 AI approaches details with example screenshots

1. Rule-Based AI Approach:

The Sinhala spell and grammar checker uses rule-based AI to correct spelling and grammar errors in Sinhala text. It leverages a preprocessed dictionary, phonetic mapping, stemming, and suffix matching to identify spelling mistakes and applies similarity metrics for corrections. The grammar checker ensures proper subject-verb agreement and word order using a POS tagger, offering accurate corrections for spelling and grammar.

```
# Input Sentences
input_sentences = [
    "අපි කෑඩ යයි.අපි උදේට කෑඩන කනි.",
    "මම අද ගදර යයි.අපි හෙට ගමට යම.",
    "අපි හෙට නුවර යන්නෙමි.මම දැන් වැට්ට යම.",
    "මම ඔයට අදරය කරයි.මම ඔයාගේ ගෙදර යයි.",
    "අපි උත්සාහයෙන් වැඩ කරමි.අපි කොහොමහර් දිනයි."
]
```

Processing Input paragraphs...

Processing: අපි කෑඩ යයි.අපි උදේට කෑඩන කනි.
Text after Spell Checking: අපි කෑඩ යයි.අපි උදේට කෑඩන් කනි.
Running Grammar Checker...

Processing: මම අද ගදර යයි.අපි හෙට ගමට යම.
Text after Spell Checking: මම අද ගෙදර යයි.අපි හෙට ගමට යම.
Running Grammar Checker...

Processing: අපි හෙට නුවර යන්නෙමි.මම දැන් වැට්ට යම.
Text after Spell Checking: අපි හෙට නුවර යන්නෙමි.මම දැන් වැට්ට යම.
Running Grammar Checker...

Processing: මම ඔයට අදරය කරයි.මම ඔයාගේ ගෙදර යයි.
Text after Spell Checking: මම ඔයට ආදරය කරයි.මම ඔයාගේ ගෙදර යයි.
Running Grammar Checker...

Processing: අපි උත්සාහයෙන් වැඩ කරමි.අපි කොහොමහර් දිනයි.
Text after Spell Checking: අපි උත්සාහයෙන් වැඩ කරමි.අපි කොහොමහර් දිනය
Running Grammar Checker...

```
Final Output: Grammar Check Results

Results for paragraph 1:
Grammar Check Result: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'මි' when subject is 'අපි'
Suggested correction: අපි කමඩි යමු
Grammar Check Result: The sentence is grammatically correct.

Results for paragraph 2:
Grammar Check Result: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'මි' when subject is 'මම'
Suggested correction: මම ඉදි ගෙදර යමි
Grammar Check Result: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'මි' when subject is 'අපි'
Suggested correction: අපි හෙට ගමට යමු

Results for paragraph 3:
Grammar Check Result: Grammar error: Verb 'යන්නෙමි' (base: යන්නෙ, affix: මි) should end with 'මි' when subject is 'අපි'
Suggested correction: අපි හෙට නුවර යන්නෙමු
Grammar Check Result: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'මි' when subject is 'මම'
Suggested correction: මම දැන් වැට්ට යමි

Results for paragraph 4:
Grammar Check Result: Grammar error: Verb 'කරයි' (base: කර, affix: යි) should end with 'මි' when subject is 'මම'
Suggested correction: මම ඕයාට ආදරය කරමි
Grammar Check Result: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'මි' when subject is 'මම'
Suggested correction: මම ඕයාගේ ගෙදර යමි

Results for paragraph 5:
Grammar Check Result: Grammar error: Verb 'කරමි' (base: කර, affix: මි) should end with 'මි' when subject is 'අපි'
Suggested correction: අපි උත්සාහයෙන් වැඩ කරමු
Grammar Check Result: Grammar error: Verb 'දිනයි' (base: දින, affix: යි) should end with 'මි' when subject is 'අපි'
Suggested correction: අපි කොහොමහර් දිනමු
```

2. Retrieval-Augmented Generation (RAG) approach:

The Sinhala grammar and spell checker uses a RAG approach with LangChain and Google's Gemini model to detect and correct errors in Sinhala sentences. It processes grammar rules from a JSON file, splits them into chunks for embedding, and stores them in a FAISS vector store for efficient retrieval. Upon receiving an input sentence, the system retrieves relevant grammar rules, constructs a prompt, and uses Gemini to generate corrections and explanations. The combination of LangChain, FAISS, and Gemini allows for accurate and fast error detection and correction.

```
Response:
Okay, let's break down the errors and correct the sentences.

**Sentence 1: අපි හෙට නුවර යන්නෙමි.**
* **Error:** The verb ending "මි" is incorrect for the subject "අපි" (we). According to the provided context, "අපි" in the present tense should use endings "මු" or "මු".
* **Corrected Sentence:** අපි හෙට නුවර යමු. (We will go to Kandy tomorrow.)
* **Explanation:** The correct verb ending for "අපි" in the present tense is "මු". The sentence is now grammatically correct.

**Sentence 2: මම දැන් වැට්ට යමි.**
* **Error:** The verb ending "මි" is incorrect for the subject "මම" (I). According to the provided context, "මම" in the present tense should use endings "මි" or "මි".
* **Corrected Sentence:** මම දැන් වැට්ට යමි. (I am going to the lake now.)
* **Explanation:** The correct verb ending for "මම" in the present tense is "මි". The sentence is now grammatically correct.

**Summary of Corrections and Errors:**
* **Original Sentence 1:** අපි හෙට නුවර යන්නෙමි.
  * **Error:** Incorrect verb ending "මි" for subject "අපි".
  * **Corrected Sentence 1:** අපි හෙට නුවර යමු.
* **Original Sentence 2:** මම දැන් වැට්ට යමි.
  * **Error:** Incorrect verb ending "මි" for subject "මම".
  * **Corrected Sentence 2:** මම දැන් වැට්ට යමි.

The provided context helped us identify the correct verb endings based on the subject and tense. The dative case marker "ට" was used correctly in the second sentence, so no changes were
```

Response:
****Corrected Sentences:****

- * මම ඉද ගෙදර යමි.
- * අපි හෙට ගමට යමු.

****Error Explanations:****

****Sentence 1: "මම ඉද ගෙදර යමි."****

- * ****Error:**** Subject-verb agreement. The subject "මම" (I) is first person singular. The verb ending "යි" is used for third person singular (he/she/it) in the present tense.
- * ****Correction:**** The correct verb ending for first person singular present tense is "මි". Therefore, "යමි" should be corrected to "යමි".

****Sentence 2: "අපි හෙට ගමට යමි."****

- * ****Error:**** Subject-verb agreement. The subject "අපි" (we) is first person plural. The verb ending "මි" is not the correct ending for first person plural present tense.
- * ****Correction:**** The correct verb ending for first person plural present tense is "මු". Therefore, "යමි" should be corrected to "යමු".

****Summary of Errors:****

Both sentences had errors related to subject-verb agreement. The verb endings did not match the person and number of the subject. The first sentence used a third-person singular ending

Response:
****Corrected Sentences:****

- * අපි කෑමේ **යමු**.
- * අපි උදේට කෑමේ**ත්** කමු.

****Explanation of Errors:****

- **Subject-Verb Agreement Error in "අපි කෑමේ යමි":****
 - * The subject "අපි" (we) is first-person plural.
 - * The verb "යමි" (goes) uses the third-person singular/plural ending "යි".
 - * According to the provided context, first-person plural requires the "මු" ending.
 - * Therefore, the correct verb ending for "අපි" is "යමු".
 - * The error type is "subject_verb_agreement" and the pattern is "subject + incorrect_verb_ending".
- **Case Marking Error in "අපි උදේට කෑමේ කමු":****
 - * The sentence "අපි උදේට කෑමේ කමු" means "We eat at the shop in the morning".
 - * The word "කෑමේ" (shop) is used as the location where the action of eating takes place.
 - * The provided context indicates that the dative case marker "ට" is used for indirect object marking. However, the location is not an indirect object.
 - * The correct case marker for location is the ablative case marker "ත්".
 - * Therefore, the correct form is "කෑමේත්" (from the shop).
 - * The error pattern is using "ත" instead of "ත්".

Response:
****Corrected Sentences:****

- * අපි උත්කරයෙන් වැඩ කරමු.
- * අපි කොහොමයේ දිනමු.

****Explanation of Errors:****

- **අපි උත්කරයෙන් වැඩ කරමි****
 - * ****Error:**** The verb ending "මි" is used with the subject "මම" (I) in the present tense. The subject here is "අපි" (we).
 - * ****Correction:**** The correct verb ending for "අපි" in the present tense is "මු" or "මු". Therefore, "කරමි" should be "කරමු".
- **අපි කොහොමයේ දිනමි****
 - * ****Error:**** The verb ending "මි" is not used with the subject "අපි" (we) in the present tense.
 - * ****Correction:**** The correct verb ending for "අපි" in the present tense is "මු" or "මු". Therefore, "දිනමි" should be "දිනමු".

****In summary:****

The errors in the original sentences were due to incorrect verb endings being used with the subject "අපි". The correct verb endings for "අපි" in the present tense are "මු" or "මු".

Response:
****Corrected Sentences:****

- * මම මයාට ආදරෙයි. (Mama oyāta ādarei.)
- * මම මයාගේ ගෙදර යමි. (Mama oyāgē gedara yamī.)

****Error Explanations:****

- **මම මයාට ආදරෙයි කරයි (Mama oyāta ādaraya karayī) - Subject-Verb Agreement Error:****
 - * ****Error:**** The verb ending "කරයි" (karayī) is used for the third-person singular subject ("ඔහු/ඇය" - he/she). However, the subject of the sentence is "මම" (mama), which is first-person singular.
 - * ****Correction:**** The correct verb ending for the first-person singular in the present tense is "මි" (mi) or "යි" (yi). The verb "ආදරෙයි කරනවා" (ādaraya karanawa) is often short.
 - * ****Explanation:**** The error falls under the "subject_verb_agreement" error type, specifically the pattern "subject + incorrect_verb_ending". The provided context also indicates
- **මම මයාගේ ගෙදර යමි (Mama oyāgē gedara yayi) - Subject-Verb Agreement Error:****
 - * ****Error:**** Similar to the first sentence, the verb ending "යමි" (yayi) is used for the third-person singular subject. The subject is "මම" (mama), which requires a first-person verb ending.
 - * ****Correction:**** The correct verb ending for the first-person singular in the present tense is "මි" (mi). Therefore, the correct form is "මම මයාගේ ගෙදර යමි" (Mama oyāgē gedara yamī).
 - * ****Explanation:**** This error also falls under the "subject_verb_agreement" error type, specifically the pattern "subject + incorrect_verb_ending". The provided context also indi

****In summary:****

Both sentences had subject-verb agreement errors. The verb endings were incorrectly conjugated for the third-person singular instead of the first-person singular. The corrected sentences

3. Fine-tuned LLM approach with the XLM-RoBERTa model:

The Sinhala grammar and spell checker uses a fine-tuned XLM-RoBERTa model for classifying sentences as grammatically correct or incorrect. It processes sentences through tokenization, fine-tunes the model with the AdamW optimizer, and evaluates grammar issues like subject-verb agreement and sentence structure. The trained model is optimized and ready for deployment to detect and suggest corrections for errors in Sinhala sentences.

Sentence: අපි කෑම යයි
Status: Incorrect
Suggestion: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'මු' when subject is 'අපි'
Suggested correction: අපි කෑම යමු

Sentence: අපි උදෙසා කෑමක කමු
Status: Incorrect
Suggestion: The sentence is grammatically correct.

Sentence: මම අද ගෙදර යයි
Status: Incorrect
Suggestion: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'මි' when subject is 'මම'
Suggested correction: මම අද ගෙදර යමි

Sentence: අපි හෙට ගමට යමි
Status: Incorrect
Suggestion: Grammar error: Verb 'යමි' (base: ය, affix: මි) should end with 'මු' when subject is 'අපි'
Suggested correction: අපි හෙට ගමට යමු

Sentence: අපි හෙට නුවර යන්නෙමි
Status: Incorrect
Suggestion: Grammar error: Verb 'යන්නෙමි' (base: යන්නෙ, affix: මි) should end with 'මු' when subject is 'අපි'
Suggested correction: අපි හෙට නුවර යන්නෙමු

Sentence: මම දැන් වැඩට යමි
Status: Incorrect
Suggestion: Grammar error: Verb 'යමි' (base: ය, affix: මි) should end with 'මි' when subject is 'මම'
Suggested correction: මම දැන් වැඩට යමි

Activate W

Sentence: මම මියාට අදරය කරයි
Status: Incorrect
Suggestion: Grammar error: Verb 'කරයි' (base: කර, affix: යි) should end with 'මි' when subject is 'මම'
Suggested correction: මම මියාට අදරය කරමි

Sentence: මම මියාගේ ගෙදර යයි
Status: Incorrect
Suggestion: Grammar error: Verb 'යයි' (base: ය, affix: යි) should end with 'ම' when subject is 'මම'
Suggested correction: මම මියාගේ ගෙදර යමි

Sentence: අපි උත්සාහයෙන් වැඩ කරමි
Status: Incorrect
Suggestion: Grammar error: Verb 'කරමි' (base: කර, affix: මි) should end with 'මු' when subject is 'අපි'
Suggested correction: අපි උත්සාහයෙන් වැඩ කරමු

Sentence: අපි කොහොමෙහි දිනයි
Status: Incorrect
Suggestion: Grammar error: Verb 'දිනයි' (base: දින, affix: යි) should end with 'මු' when subject is 'අපි'
Suggested correction: අපි කොහොමෙහි දිනමු

c. Selected approach and why

Selected approach: Rule based approach

The Rule-Based AI Approach was chosen due to its high accuracy of 90%, outperforming the LLM model at 70% and the RAG model at 60%. By using predefined rules, it ensures precise handling of Sinhala grammar and sentence structure, making it ideal for tasks requiring accuracy. The limited availability of large Sinhala datasets negatively impacted the performance of the LLM and RAG models, restricting their ability to generalize and learn diverse language patterns.

Although the LLM and RAG models offer greater flexibility, they struggled with issues like subject-verb agreement and sentence structure errors. For instance, the LLM model failed to produce correct sentences such as "අපි කෑමක කමු" and "මම අද ගෙදර යමි." The lack of sufficient training data exacerbated these problems, leading to lower reliability compared to the Rule-Based model. As a result, the Rule-Based approach emerged as the more dependable choice for accurate Sinhala language processing.

4. Conclusion:

In conclusion, the Rule-Based Model achieved the highest accuracy of 90%, outperforming both the LLM and RAG models, which scored 70% and 60%, respectively. The Rule-Based Model's precision is due to its reliance on predefined rules, ensuring consistent grammar and sentence structure. In contrast, while the LLM and RAG models offer flexibility, they struggled with issues such as incorrect subject-verb agreements and inconsistent sentence structures, leading to errors in output. These results highlight the trade-off between flexibility and precision, suggesting that rule-based systems are ideal for tasks requiring strict language rules, while LLM and RAG models can be useful for more dynamic language processing with further refinement.

```
correct_outputs = [  
    "අපි කඩේට යමු",  
    "අපි උදේට කඩෙන් කමු",  
    "මම අද ගෙදර යමි",  
    "අපි හෙට ගමට යමු",  
    "අපි හෙට නුවර යන්නෙමු",  
    "මම දැන් වැවට යමි",  
    "මම ඔයාට ආදරය කරමි",  
    "මම ඔයාගේ ගෙදර යමි",  
    "අපි උත්සාහයෙන් වැඩ කරමු",  
    "අපි කොහොමහරි දිනමු"  
]
```

Accuracy Scores:
LLM Model: 70.00%
RAG Model: 60.00%
Rule-based Model: 90.00%

Errors in LLM Model:

Line 1:
Expected: අපි කඩේට යමු
Got : අපි කඩේ යමු

Line 2:
Expected: අපි උදේට කඩෙන් කමු
Got : අපි උදේට කඩෙන් කමු

Line 3:
Expected: මම අද ගෙදර යමි
Got : මම අද ගෙදර යමි

Errors in RAG Model:

Line 1:
Expected: අපි කඩේට යමු
Got : අපි කඩේ යමු

Line 2:
Expected: අපි උදේට කඩෙන් කමු
Got : අපි උදේට කඩෙන් කමු

Line 5:
Expected: අපි හෙට නුවර යන්නෙමු
Got : අපි හෙට නුවර යමු

Line 7:
Expected: මම ඔයාට අදරය කරමි
Got : මම ඔයාට අදරය කරනවා

Errors in Rule-based Model:

Line 1:
Expected: අපි කඩේට යමු
Got : අපි කඩේ යමු

GitHub link: <https://github.com/SalindaHulangamuwa/Spelling-corrector-and-grammar-checker-for-Sinhala>