

# Team Binary Brigade

Team Leader - Phyo Kyi  
Team Member - Thant Thiri Maung

# **Distinguishing Burmese Male and Female Names**

# Team Binary Brigade

Our team



**Phyo Kyi**

Team leader

0

1

**Thant Thiri Maung**

Team Member



# About Project



## Our Purpose

To classify Myanmar Names by sex (male or female or both).



## Impact

Useful for Myanmar Name Entity Recognition



## Motivation

This can be great impact in Myanmar NER and this become our motivation

# AI Ethics

- **Fairness**
  - Male and female name size nearly equal.
  - Not only on emphasize the Burmese ethnic group, but other also
  - No left for the names which can be both
- **Reliability and Safety - 93.08176100628931% accuracy score**
- **Privacy and Security - We only take name data**



# AI Ethics

- **Inclusiveness**
  - **Work on some ethnic names (Shan , Kayin), plan to include others**
  - **Replacing similar words help user wrong input to correct**
- **Transparency**
  - **can use Myanmar and English Language**
  - **Most common burmese names only**
  - **Unicode Only**
- **Accountability**
  - **Public Github Repo**
  - **Google Form for New Names**

# Data Collection

- **From Social Media, from GitHub, other sources like articles**

[GitHub - L16H7/Myanmar\\_Names: Open Source List of Myanmar\(Burmese\) Names for Male & Female Names](#)

- **Total Data - 6357, Male = 3191, Female = 3172**
- **English > Burmese Name changing**

<https://docs.google.com/spreadsheets/d/1FRWeY1QMEsyDTOGjk4Jj5hGXQekQb0-uURDbuAKe3Cs/edit#gid=1461949965>

# Data Encoding

Give the label male as 0 and female as 1.



# Data for Rule based Methods

- Special female two words name (eg, သီရိ, သဉ္ဇာ, သိင်္ဂီ ) > 105
- Special male two words name (eg, သီဟ, သူရ) > 5
- Both words (eg, သန်းဆွေ, သန်းမြင့်အောင်) > 91

# Data Preparation

## Tokenization

- By Segment ['ခိုင်', 'သန္တာ', 'ထွန်း']
- By Character ['ခ', 'ိ', 'င်', 'သ', 'န', 'တ', 'ာ', 'ထွ', 'န်း']
- By Syllable Tokenization ['ခိုင်', 'သ', 'န္တာ', 'ထွန်း']
- By Multilingual Semi Syllable Break ['ခိ', 'င်', 'သ', 'န', 'တာ', 'ထွ', 'န်း']

Sources

<https://github.com/swanhtet1992/ReSegment/blob/master/resegment.py>

<https://github.com/SaPhyoThuHtet/nlp-tool/blob/main/utilities.py>

# Replacing similar words

```
def clean_text(text):  
    text = text.replace("\u200c", "")  
    text = text.replace(" ", "")  
    text = text.replace("ချ", "၍")  
    text = text.replace("ည်", "ီ")  
    text = text.replace("ဏ်", "န်")  
    text = text.replace("မ်", "န်")  
    text = text.replace("ါ", "ာ")  
    text = text.replace("အူ", "ဦး")  
    return text
```

# Data Splitting

- **train\_test\_split**

```
max_i=0
max_score=0
for i in range(10,70):
    x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size = random_size,random_state=i)
    classifier = LogisticRegression()
    if model == "mf_logistic_regression_by_character_tokenization":
        classifier = LogisticRegression(max_iter=5000)
    classifier.fit(x_train, y_train)      #training level
    score = classifier.score(x_test, y_test)
    if score > max_score:
        max_i=i
        max_score=score
        max_classifier = classifier
    max_x_train, max_x_test, max_y_train, max_y_test = x_train, x_test, y_train, y_test
    max_y_pred = classifier.predict(max_x_test) # Predit Y Value with Test Datasets
```

# Understanding on Algorithm

Rule based Systems

Logistics Regression

Dictionary



# Logistics Regression

- Use in Binary/Categorical Classification
- Logistics Regression use Logistic function/ Sigmoid function

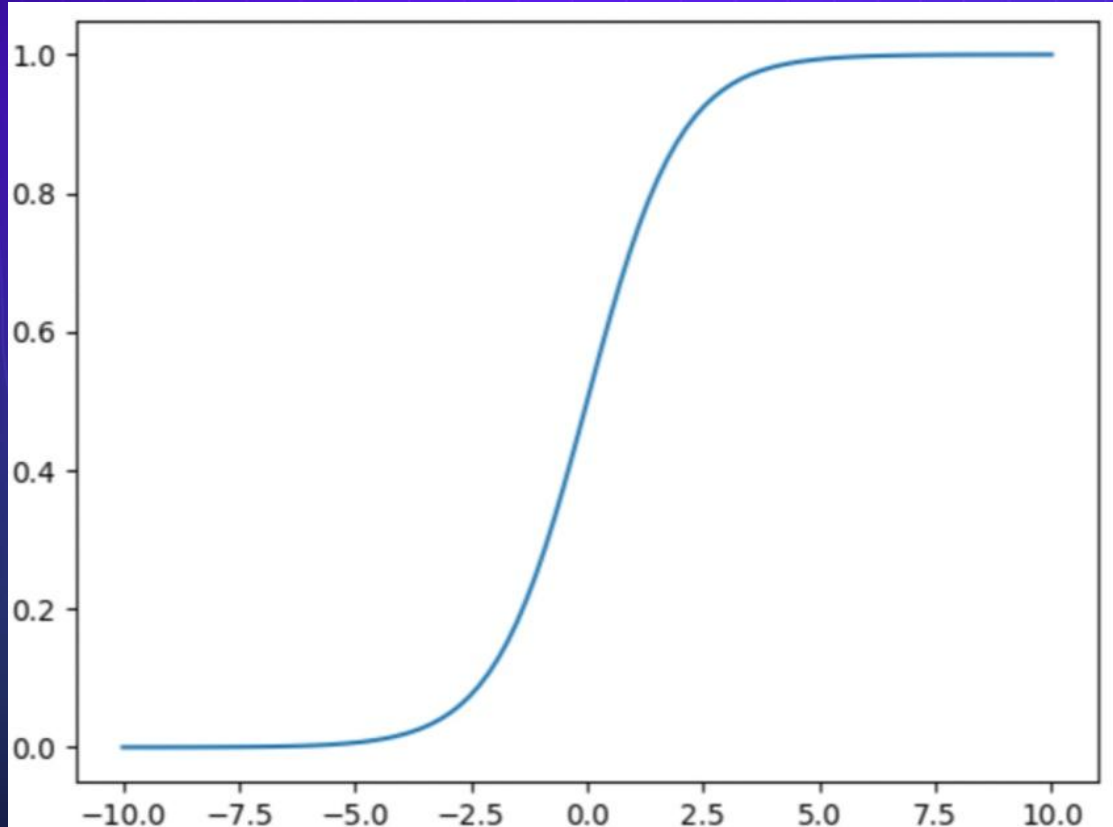
Formula

$$S(x) = \frac{1}{1 + e^{-x}}$$

$S(x)$  = sigmoid function

$e$  = Euler's number

# S-shape curve



Input values between minus and plus infinity, only values between 0 and 1 result.

Then, Sigmoid function use threshold value which can decide that a given input belongs to what type of two classes.

|  |                 |        | ခင် | စိုး | မိုး | ဝင်း | နိုင် | စု | sex |
|--|-----------------|--------|-----|------|------|------|-------|----|-----|
|  | ခင်စိုးမိုးဝင်း | female | 1   | 1    | 1    | 1    | 0     | 0  | 1   |
|  | စိုးမိုးနိုင်   | male   | 0   | 1    | 1    | 0    | 1     | 0  | 0   |
|  | မိုးမိုးစု      | female | 0   | 0    | 2    | 0    | 0     | 1  | 1   |

| Syllable 1 |          | Syllable 2 |          | Syllable 3 |          | sex |
|------------|----------|------------|----------|------------|----------|-----|
| Male %     | Female % | Male %     | Female % | Male %     | Female % |     |
| 13         | 100      | 30         | 50       | 20         | 10       | 1   |
| 100        | 20       | 30         | 20       | 0          | 0        | 0   |
| 40         | 100      | 40         | 100      | 10         | 100      | 1   |

# Model and Results

**By Character Tokenization - 86.9496855345912%**

**By Multilingual Semi Syllable - 93.08176100628931%**

**By Segment - 92.45283018867924%**

**By Syllable Tokenization - 93.08176100628931%**

# Experiment

- **Logistics Regression with Tokenization: Character, Syllable, Segment, Multi Syllable**
- **Rule Base (Special Word and Leading Two Words) and Dictionary**
- **Logistics Regression weight value on Syllable**
- **Tensorflow Neural Network**

```
Clean Data is finished.
Word Columns by Segments is finish
Word Columns by Syllable Tokenization is finished.
Word Columns by Multilingual Semi Syllable Break is finished.
Word Columns by Character Tokenization is finished.
MF leading exclusion list is finished.
=====
By Segment
=====
Total is 6358. Male is 3189. Female is 3169
Model accuracy score of random state 11 is : 92.45283818867924
[[295 16]
 [ 32 293]]
=====
By Syllable Tokenization
=====
Total is 6358. Male is 3189. Female is 3169
Model accuracy score of random state 65 is : 93.08176108628931
[[388 24]
 [ 28 284]]
=====
By Multilingual Semi Syllable Break
=====
Total is 6358. Male is 3189. Female is 3169
Model accuracy score of random state 54 is : 93.08176108628931
[[304 17]
 [ 27 288]]
=====
By Character Tokenization
=====
Total is 6358. Male is 3189. Female is 3169
Model accuracy score of random state 54 is : 86.9496855345912
[[276 45]
 [ 38 277]]
```

```
=====test and generate for model segment=====
Right is 6056
Wrong is 304
Percentage is 95.22012578616352

Wrong Method Count
{0: 0, 1: 0, 2: 0, 3: 0, 4: 0, 5: 304, 6: 0}
Right Method Count
{0: 0, 1: 1152, 2: 31, 3: 1008, 4: 2752, 5: 1026, 6: 87}
=====
=====test and generate for model syllable=====
Right is 6100
Wrong is 260
Percentage is 95.9119496855346

Wrong Method Count
{0: 0, 1: 0, 2: 0, 3: 0, 4: 0, 5: 260, 6: 0}
Right Method Count
{0: 0, 1: 1152, 2: 31, 3: 1008, 4: 2752, 5: 1070, 6: 87}
=====
=====test and generate for model multi=====
Right is 6246
Wrong is 114
Percentage is 98.20754716981132

Wrong Method Count
{0: 0, 1: 0, 2: 0, 3: 0, 4: 0, 5: 114, 6: 0}
Right Method Count
{0: 0, 1: 1152, 2: 31, 3: 1008, 4: 2752, 5: 1216, 6: 87}
=====
=====test and generate for model character=====
Right is 5494
Wrong is 866
Percentage is 86.38364779874213

Wrong Method Count
{0: 0, 1: 0, 2: 0, 3: 0, 4: 0, 5: 866, 6: 0}
Right Method Count
{0: 0, 1: 1152, 2: 31, 3: 1008, 4: 2752, 5: 464, 6: 87}
=====
```

Speaker - 0



# Data for Rule based Methods

- M6 : Both words (eg, သန်းဆွေ, သန်းမြင့်အောင်) > 91
- M0 : Special List
- M1 : Special female two words name (eg, သီရိ, သဉ္ဇာ, သိင်္ဂီ ) > 105
- M2 : Special male two words name (eg, သီဟ, သူရ) > 5
- M3 : Leading words for female
- M4 : Leading words for male
- M5 : from Model

# Streamlit Demo

<https://myanmar-names-male-female.streamlit.app>

## Github

<https://github.com/binarybrigade/myanmar-name-mf-classification>

**Thank You**

# Any Questions?

