

**ĐẠI HỌC QUỐC GIA TP. HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN**



BÀI TẬP MÔN CÁC KỸ THUẬT HỌC SÂU VÀ ỨNG DỤNG

KHOA: KHOA HỌC MÁY TÍNH

HOMEWORK: TEXT GENERATION WITH RNN FOR VIETNAMESE

Giảng viên: Nguyễn Duy Khánh

Nhóm thực hiện:

- 1. Trương Văn Khải– 21520274**

Đề: Làm Text Generation với RNN cho dữ liệu Tiếng Việt tự thu thập, tối thiểu 1000 ký tự

Bài làm:

- **RNN cho cấp độ âm tiết (characters)**

Source: [Tạo văn bản bằng RNN | Text | TensorFlow](#)

1. Dataset:

- Lời nhạc 3 bài hát của ca sĩ Đạt G: Thêm Bao Nhiêu Lâu + Khó Về Nụ Cười + Bánh Mì Không
- Tiền xử lý dataset:
 - o Độ dài: 6354 characters
 - o Vocabulary size: 107 unique characters

2. Mô tả kiến trúc Model:

- Embedding Input: 107 embedding vectros, mỗi vector có chiều = embedding_dim = 256
- Rnn_units = 1024
- Output: Dense layer : 97

3. Huấn luyện:

- batch_size = 2 (Với dữ liệu ít, để batch_size cao quá mô hình dự đoán ra rất tệ)
- EPOCHS = 50
- Optimizer = 'adam'
- loss = SparseCategoricalCrossentropy

4. Kết quả sinh text:

Số lượng từ cần sinh là: 500

- Trường hợp 1: Từ bắt đầu là “**á**”
- Kết quả: **ánh mắt anh buồn Hoạ nụ cười thêm trên nét môi Nhưng sao chẳng thương có thể em sống cho em Quá dễ để biết tên nhau nhưng đâu có thể nhìn thấy nhau đau ời hiểu thấu Vẫn hoạ thêm chiếc môi cười tiếp theo Nhưng đau thấu trời Khóc thật nhiều Ngồi khóc thật nhiều Khóc cho đời phong ba lăm đau mà chông thế này. ... Anh mơ một mai thức giấc chẳng cơn đau nào sẽ ghé qua Em mong ngày mai có nắng sưởi ấm con tim buồn đau tối qua Anh mong dù cho vấp ngã có lăm phong ba bình yên cho em nhưng cũng có thể em số**

```
start = time.time()
status = None
next_char = tf.constant(['á'])
result = [next_char]

for n in range(500):
    next_char, states = one_step_model.generate_one_step(next_char, states=states)
    result.append(next_char)

result = tf.strings.join(result)
end = time.time()
print(result[0].numpy().decode('utf-8'), '\n' + ' ' * 80)

print('Length text generation: ', len(result[0].numpy().decode('utf-8')))
print('Value time: ', end - start)
```

á ánh mắt anh buồn Hoạ nụ cười thêm trên nét môi Nhưng sao chẳng thương có thể em sống cho em Quá dễ để biết tên nhau nhưng đau có thể nhìn thấy nhau đau ời hiểu thấu Vẫn hoạ thêm chiếc môi cười tiếp theo M

Length text generation: 501

Run time: 1.259234286089111

- Trường hợp 2: Từ bắt đầu là “**.**”
- Kết quả: **. D em xem mì côi đưng cho em Quá dễ để biết tên nhau nhưng đâu có thể nhìn thấy nhau đau Chiều hoàng hôn buông xuống phía tây nghen ngào Uống phai oh Nỗi đau dày vò nỗi đau trời cao có thấu, có biết cơn đau mùi hương thế nào. ... Biết đâu những cơn say đau đời này Phải khóc cho đến hôm nay Hẹn ước chi thể em ơi giờ phải vậy Gửi gió nổi**

nhớ bay bay Nhìn hoàng hôn kia còn mang em đi xa anh
thêm bao nhiêu lâu nữa sẽ trả về Đừng có như vậy, dặn lòng
đừng khóc như vậy Mà đời đâu như anh mơ, đâu như anh

```
start = time.time()
states = None
next_char = tf.constant([' '])
result = [next_char]

for n in range(500):
    next_char, states = rnn_step_model.generate_one_step(next_char, states=states)
    result.append(next_char)

result = tf.strings.join(result)
end = time.time()
print(result.numpy().decode('utf-8'), '\n\n' + '*'*80)

print("length text generation: ", len(result[0].numpy().decode('utf-8')))
print('\nRun time:', end - start)
```

Dm xem nì cái đng cho em quá dễ để biết tên nhau nhưng đầu có thể nhìn thấy nhau đau chiều hoàng hôn bóng xuống phía tây ghen ngào từng phai oh nỗi đau day và nỗi đau trời cao có thấu, có biết em d

length text generation: 501

Run time: 2.0114852985273458

• RNN cho cấp độ từ

1. Dataset

- Lời nhạc 6 bài hát của ca sĩ Đạt G: Thêm Bao Nhiêu Lâu + Khó Về Nụ Cười + Bánh Mì Không Ngày Mai Em Đi Mất + Còn Buồn Không Em + Anh Tự Do Nhưng Cô Đơn
- Tiền xử lý dataset:
 - o Độ dài: 493 từ (11534 characters)
 - o Vocabulary size: 493 unique characters

2. Mô tả kiến trúc Model

- Embedding Input: 493 embedding vectros, mỗi vector có chiều = embedding_dim = 256
- Một lớp 1 Bidirectional với cell LSTM 150 node
- Một lớp LSTM với 100 node
- Output: Dense layer : 493

3. Huấn luyện

- EPOCHS = 150
- Optimizer = 'adam'
- loss = categorical_crossentropy

4. Kết quả sinh text:

- Số lượng từ cần sinh là: 200
- Từ bắt đầu “em ỏn không”
- Kết quả: **em ỏn không hiểu thấu em như nắng mai trong đời cho em đau ai ngờ để phải nhớ có có khóc xôi ai anh thêm bao nhiêu lâu anh có thêm rơi sẽ hẹn rong chơi và côỉ rồi nhạc vang lên xe tranh và côỉ lắỉ cứ rơi nhưng có yêu ai sao sao đêm nay trong xa là ai anh sao anh khóc cho qua thấu đừng khóc có có thấu đau thế ai ai cười lắỉ vương em có cơn khóc có có xưa còn nói nhớ có vào ủa về và cứ rơi anh tối chẳng khi tổ như mưa rơi đâu em hiểu anh dành hơn nữa thứ thứ đó em ở nhớ em có sẽ lắỉ ở phần anh ỏi lại thể anh sẽ chờ em ỏi lại cho anh biết vương và vắỉ khô rồi kêu vắỉ mà vắỉ phần cũng không biết bầu trời kia lại anh dành thứ chậm em sao anh cũng nhắm mắt em cũng khóc còn vào kiểỉm ắỉ rơi còn có cơn nhưng có nhau ắỉ xôi lắỉ này đời phiểỉn thắỉng trắỉm ắỉ về rơi nhưng mưa rơi nếu thấu em mưa rơi em yêu như thế mà**

5. Predict the next N words.

```
[42] test_seq = 'em ớn không'
```

```
next_words = 200
```

```
for _ in range(next_words):
    token_list = tokenizer.texts_to_sequences([test_seq])[0]
    token_list = pad_sequences([token_list], maxlen=max_sequence_len-1, padding='pre')
    predicted = model.predict(token_list, verbose=0)

    output_word = ""
    predicted_id = np.argmax(predicted)

    if predicted_id in tokenizer.index_to_word:
        output_word = tokenizer.index_to_word[predicted_id]

        if output_word == '<end>':
            break
        test_seq += " " + output_word
    else:
        break

print(test_seq)
```

em ớn không hiểu thấu em như nắng mai trong đời cho em đau ai ngờ đó phải nhớ có cô khóc với ai anh thêm bao nhiêu lâu anh có thêm rồi sẽ hẹn rong chơi và rồi rồi nhấc vang lên xe tr