

# Bài 2 - Lý thuyết về mạng LSTM part 2

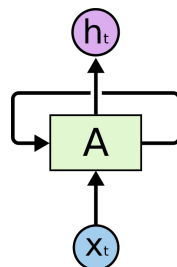
22 Apr 2019 - phamdinhhkhanh

## Menu

- 1. Mạng nơ ron truy hồi (RNN - Recurrent Neural Network)
- 2. Hạn chế của mạng nơ ron truy hồi
- 3. Mạng trí nhớ ngắn hạn định hướng dài hạn (LSTM - Long short term memory)
- 4. Ý tưởng đằng sau LSTM
- 5. Thứ tự các bước của LSTM
- 6. Các biến thể của LSTM
- 7. Kết luận
- 8. Thực hành mô hình sinh từ tự động
  - 8.1 Xây dựng mô hình trên level kí tự
  - 8.2. Mô hình LSTM
  - 8.3. Kiến trúc BiLSTM.
- 9. Tài liệu

## 1. Mạng nơ ron truy hồi (RNN - Recurrent Neural Network)

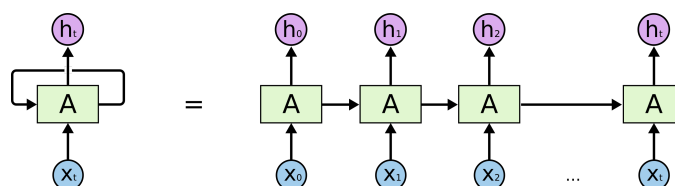
Trong lý thuyết về ngôn ngữ, ngữ nghĩa của một câu được tạo thành từ mối liên kết của những từ trong câu theo một cấu trúc ngữ pháp. Nếu xét từng từ một đứng riêng lẻ ta không thể hiểu được nội dung của toàn bộ câu, nhưng dựa trên những từ xung quanh ta có thể hiểu được trọn vẹn một câu nói. Như vậy cần phải có một kiến trúc đặc biệt hơn cho các mạng nơ ron biểu diễn ngôn ngữ nhằm mục đích liên kết các từ liên trước với các từ ở hiện tại để tạo ra mối liên hệ sâu chuỗi. Mạng nơ ron truy hồi đã được thiết kế đặc biệt để giải quyết yêu cầu này:



**Hình 1: Mạng nơ ron truy hồi với vòng lặp**

Hình trên biểu diễn kiến trúc của một mạng nơ ron truy hồi. Trong kiến trúc này mạng nơ ron sử dụng một đầu vào là một véc tơ  $x_t$  và trả ra đầu ra là một giá trị ẩn  $h_t$ . Đầu vào được đầu vào với một thân mạng nơ ron  $A$  có tính chất truy hồi và thân này được đầu vào tới đầu ra  $h_t$ .

Vòng lặp  $A$  ở thân mạng nơ ron là điểm mấu chốt trong nguyên lý hoạt động của mạng nơ ron truy hồi. Đây là chuỗi sao chép nhiều lần của cùng một kiến trúc nhằm cho phép các thành phần có thể kết nối liền mạch với nhau theo mô hình chuỗi. Đầu ra của vòng lặp trước chính là đầu vào của vòng lặp sau. Nếu trải phẳng thân mạng nơ ron  $A$  ta sẽ thu được một mô hình dạng:



**Hình 2: Cấu trúc trải phẳng của mạng nơ ron truy hồi**

Kiến trúc mạng nơ ron truy hồi này tỏ ra khá thành công trong các tác vụ của deep learning như: Nhận diện giọng nói (speech recognition), các mô hình ngôn ngữ, mô hình dịch, chú thích hình ảnh (image captioning),....

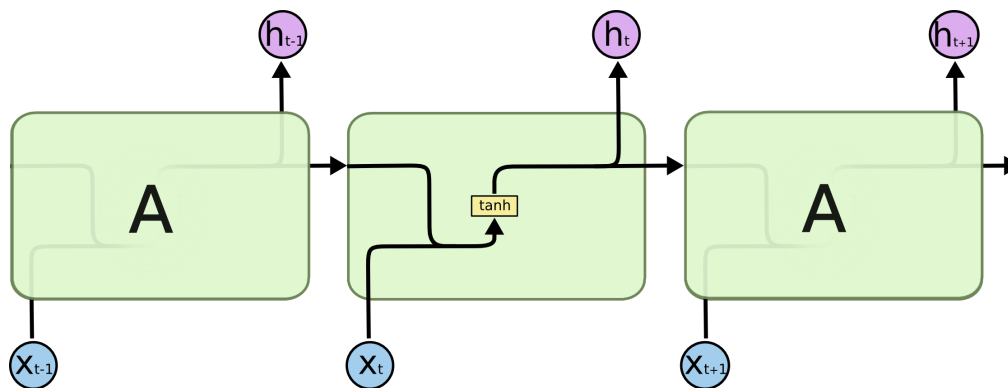
## 2. Hạn chế của mạng nơ ron truy hồi

Một trong những điểm đặc biệt của RNN đó là nó có khả năng kết nối các thông tin liên trước với nhiệm vụ hiện tại, chẳng hạn như trong câu văn: 'học sinh đang tới *trường học*'. Dường như trong một ngữ cảnh ngắn hạn, từ *trường học* có thể được dự báo ngay tức thì mà không cần thêm các thông tin từ những câu văn khác gần đó. Tuy nhiên có những tình huống đòi hỏi phải có nhiều thông tin hơn chẳng hạn như: 'hôm qua Bổng đi học nhưng không mang áo mưa. Trên đường đi học trời mưa. Cặp sách của Bổng bị *ướt*'. Chúng ta cần phải học để tìm ra từ *ướt* ở một ngữ cảnh dài hơn so với chỉ 1 câu. Tức là cần phải biết các sự kiện trước đó như *trời mưa, không mang áo mưa* để suy ra sự kiện bị *ướt*. Những sự liên kết ngữ nghĩa dài như vậy được gọi là phụ thuộc dài hạn (*long-term dependencies*). Về mặt lý thuyết mạng RNN có thể giải quyết được những sự phụ thuộc trong dài hạn. Tuy nhiên trên thực tế RNN lại cho thấy khả năng học trong dài hạn kém hơn. Để hiểu thêm lý do tại sao mạng RNN lại không có khả năng học trong dài hạn cùng đọc bài Learning Long - Term Dependencies with Gradient Descent is Difficult (<http://ai.dinfo.unifi.it/paolo/ps/tnn-94-gradient.pdf>). Một trong những nguyên nhân chính được giải thích đó là sự triệt tiêu đạo hàm của hàm cost function sẽ diễn ra khi trải qua chuỗi dài các tính toán truy hồi. Một phiên bản mới của mạng RNN là mạng LSTM ra đời nhằm khắc phục hiện tượng này nhờ một cơ chế đặc biệt.

## 3. Mạng trí nhớ ngắn hạn định hướng dài hạn (LSTM - Long short term memory)

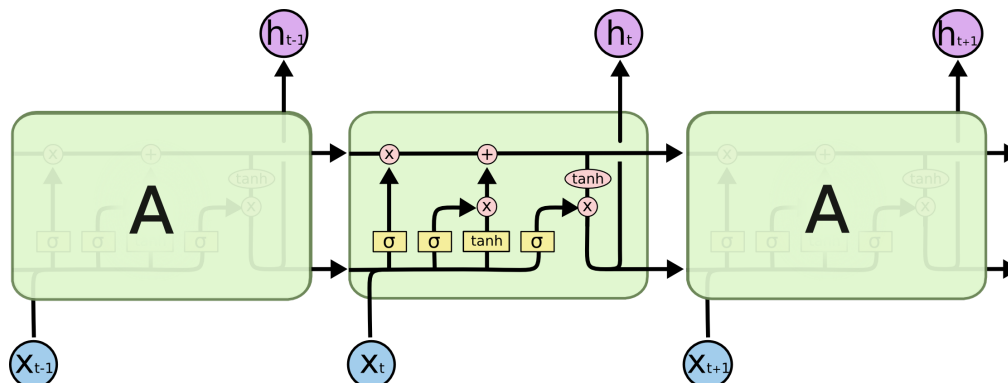
Mạng trí nhớ ngắn hạn định hướng dài hạn còn được viết tắt là LSTM là một kiến trúc đặc biệt của RNN có khả năng học được sự phụ thuộc trong dài hạn (*long-term dependencies*) được giới thiệu bởi Hochreiter & Schmidhuber (1997) (<http://www.bioinf.jku.at/publications/older/2604.pdf>). Kiến trúc này đã được phổ biến và sử dụng rộng rãi cho tới ngày nay. LSTM đã tỏ ra khắc phục được rất nhiều những hạn chế của RNN trước đây về triệt tiêu đạo hàm. Tuy nhiên cấu trúc của chúng có phần phức tạp hơn mặc dù vẫn dựa trên tư tưởng chính của RNN là sự sao chép các kiến trúc theo dạng chuỗi.

Một mạng RNN tiêu chuẩn sẽ có kiến trúc rất đơn giản chẳng hạn như đối với kiến trúc gồm một tầng ẩn là hàm tanh như bên dưới.



Hình 3: Sự lặp lại kiến trúc module trong mạng RNN chứa một tầng ẩn

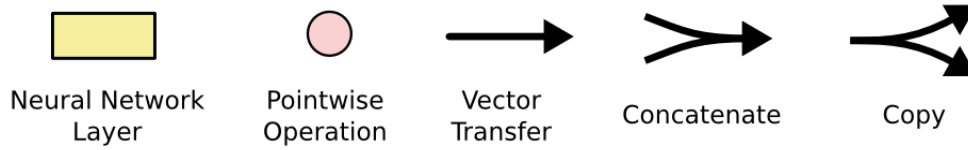
LSTM cũng có một chuỗi dạng như thế nhưng phần kiến trúc lặp lại có cấu trúc khác biệt hơn. Thay vì chỉ có một tầng đơn, chúng có tới 4 tầng ẩn (3 sigmoid và 1 tanh) tương tác với nhau theo một cấu trúc đặc biệt.



Top

#### Hình 4: Sự lặp lại kiến trúc module trong mạng LSTM chứa 4 tầng ẩn (3 sigmoid và 1 tanh) tương tác

Các kí hiệu có thể diễn giải như sau:

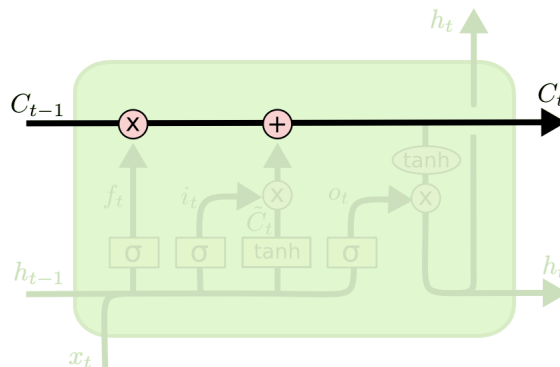


#### Hình 5: Diễn giải các kí hiệu trong đồ thị mạng nơ ron (áp dụng chung cho toàn bộ bài)

Trong sơ đồ tính toán trên, mỗi một phép tính sẽ triển khai trên một véc tơ. Trong đó hình tròn màu hồng biểu diễn một toán tử đối với véc tơ như phép cộng véc tơ, phép nhân vô hướng các véc tơ. Màu vàng thể hiện hàm activation mà mạng nơ ron sử dụng để học trong tầng ẩn, thông thường là các hàm phi tuyến sigmoid và tanh. Kí hiệu 2 đường thẳng nhập vào thể hiện phép chập kết quả trong khi kí hiệu 2 đường thẳng rẽ nhánh thể hiện cho nội dung véc tơ trước đó được sao chép để đi tới một phần khác của mạng nơ ron.

## 4. Ý tưởng đằng sau LSTM

Ý tưởng chính của LSTM là thành phần ô trạng thái (cell state) được thể hiện qua đường chạy ngang qua đỉnh đồ thị như hình vẽ bên dưới:

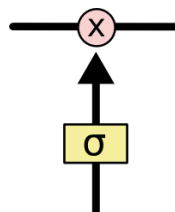


#### Hình 6: Đường đi của ô trạng thái (cell state) trong mạng LSTM

Ô trạng thái là một dạng băng chuyền chạy thẳng xuyên suốt toàn bộ chuỗi với chỉ một vài tương tác tuyến tính nhỏ giúp cho thông tin có thể truyền dọc theo đồ thị mạng nơ ron ổn định.

LSTM có khả năng xóa và thêm thông tin vào ô trạng thái và điều chỉnh các luồng thông tin này thông qua các cấu trúc gọi là cổng.

Cổng là cơ chế đặc biệt để điều chỉnh luồng thông tin đi qua. Chúng được tổng hợp bởi một tầng ẩn của hàm activation sigmoid và với một toán tử nhân như đồ thị.



#### Hình 7: Một cổng của hàm sigmoid trong LSTM

Top

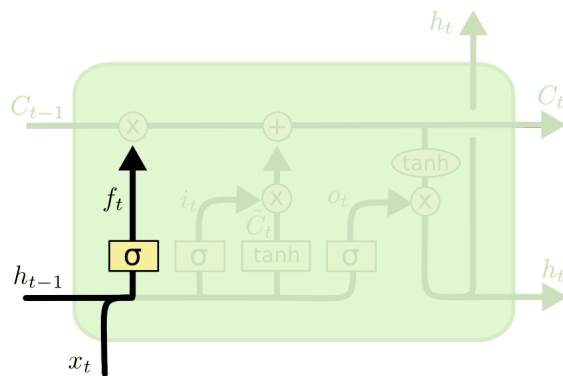
Hàm sigmoid sẽ cho đầu ra là một giá trị xác suất nằm trong khoảng từ 0 đến 1, thể hiện rằng có bao nhiêu phần thông tin sẽ đi qua cổng. Giá trị bằng 0 ngụ ý rằng không cho phép thông tin nào đi qua, giá trị bằng 1 sẽ cho toàn bộ thông tin đi qua.

Một mạng LSTM sẽ có 3 cổng có kiến trúc dạng này để bảo vệ và kiểm soát các ô trạng thái.

## 5. Thứ tự các bước của LSTM

Bước đầu tiên trong LSTM sẽ quyết định xem thông tin nào chúng ta sẽ cho phép đi qua ô trạng thái (cell state). Nó được kiểm soát bởi hàm sigmoid trong một tầng gọi là tầng quên (*forget gate layer*). Đầu tiên nó nhận đầu vào là 2 giá trị  $h_{t-1}$  và  $x_t$  và trả về một giá trị nằm trong khoảng 0 và 1 cho mỗi giá trị của ô trạng thái  $C_{t-1}$ . Nếu giá trị bằng 1 thể hiện 'giữ toàn bộ thông tin' và bằng 0 thể hiện 'bỏ qua toàn bộ chúng'.

Trở lại ví dụ về ngôn ngữ, chúng ta đang cố gắng dự báo từ tiếp theo dựa trên toàn bộ những từ trước đó. Trong những bài toán như vậy, ô trạng thái có thể bao gồm loại của chủ ngữ hiện tại, để cho đại từ ở câu tiếp theo được sử dụng chính xác. Chẳng hạn như chúng ta đang mô tả về một người bạn là con trai thì các đại từ nhân xưng ở tiếp theo phải là anh, thằng, hắn thay vì cô, con ấy (xin lỗi vì lấy ví dụ hơi thô). Tuy nhiên chủ ngữ không phải khi nào cũng cố định. Khi chúng ta nhìn thấy một chủ ngữ mới, chúng ta muốn quên đi loại của một chủ ngữ cũ. Do đó tầng quên cho phép cập nhật thông tin mới và lưu giữ giá trị của nó khi có thay đổi theo thời gian.

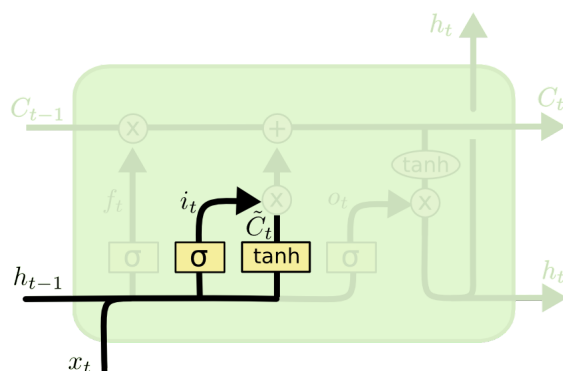


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

**Hình 8: Tầng cổng quên (*forget gate layer*)**

Bước tiếp theo chúng ta sẽ quyết định loại thông tin nào sẽ được lưu trữ trong ô trạng thái. Bước này bao gồm 2 phần. Phần đầu tiên là một tầng ẩn của hàm sigmoid được gọi là tầng cổng vào (*input gate layer*) quyết định giá trị bao nhiêu sẽ được cập nhật. Tiếp theo, tầng ẩn hàm tanh sẽ tạo ra một véc tơ của một giá trị trạng thái mới  $\tilde{C}_t$  mà có thể được thêm vào trạng thái. Tiếp theo kết hợp kết quả của 2 tầng này để tạo thành một cập nhật cho trạng thái.

Trong ví dụ của mô hình ngôn ngữ, chúng ta muốn thêm loại của một chủ ngữ mới vào ô trạng thái để thay thế phần trạng thái cũ muốn quên đi.



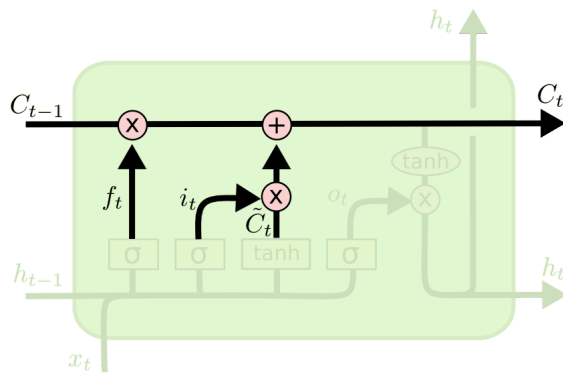
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

**Hình 9: Cập nhật giá trị cho ô trạng thái bằng cách kết hợp 2 kết quả từ tầng cổng vào và tầng ẩn hàm tanh**

Đây là thời điểm để cập nhật một ô trạng thái cũ,  $C_{t-1}$  sang một trạng thái mới  $C_t$ . Những bước trước đó đã quyết định làm cái gì, và tại bước này chỉ cần thực hiện nó.

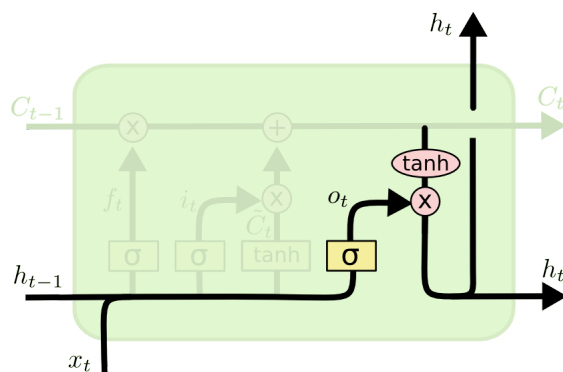
Chúng ta nhân trạng thái cũ với  $f_t$  tương ứng với việc quên những thứ quyết định được phép quên sớm. Phần tử đề cử  $i_t * \tilde{C}_t$  là một giá trị mới được tính toán tương ứng với bao nhiêu được cập nhật vào mỗi giá trị trạng thái.



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Hình 10: Ô trạng thái mới

Cuối cùng cần quyết định xem đầu ra sẽ trả về bao nhiêu. Kết quả ở đầu ra sẽ dựa trên ô trạng thái, nhưng sẽ là một phiên bản được lọc. Đầu tiên, chúng ta chạy qua một tầng sigmoid nơi quyết định phần nào của ô trạng thái sẽ ở đầu ra. Sau đó, ô trạng thái được đưa qua hàm tanh (để chuyển giá trị về khoảng -1 và 1) và nhân nó với đầu ra của một cổng sigmoid, do đó chỉ trả ra phần mà chúng ta quyết định.



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

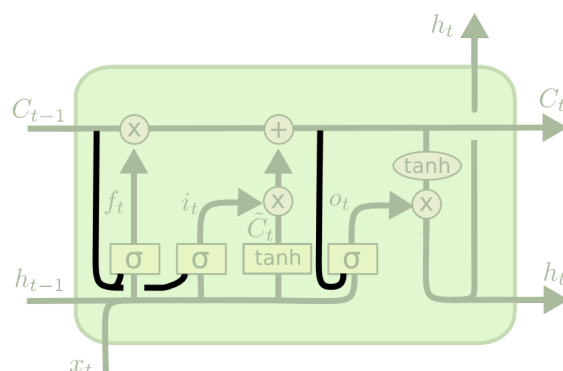
Hình 11: Điều chỉnh thông tin ở đầu ra thông qua hàm tanh

## 6. Các biến thể của LSTM

Những gì mà chúng ta vừa mô tả cho đến giờ là một mạng LSTM rất thông thường. Nhưng không phải toàn bộ LSTM đều tương tự như trên. Trên thực tế, có vẻ như hầu hết mọi bài báo liên quan đến LSTM đều sử dụng những version khác nhau đôi chút. Sự khác biệt là rất nhỏ nhưng rất đáng để đề cập một ít trong số những kiến trúc này.

Một trong những biến thể nổi tiếng nhất của LSTM được giới thiệu bởi Gers & Schmidhuber (2000)

(<ftp://ftp.idsia.ch/pub/juergen/TimeCount-IJCNN2000.pdf>) thêm một kết nối ống tiêu (peehole connection) để các cổng có thể kết nối trực tiếp đến các ô trạng thái.



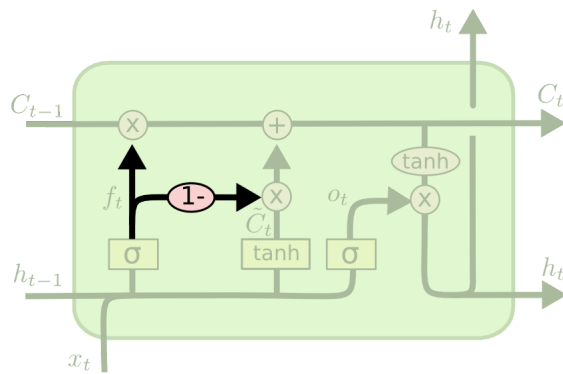
$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

Hình 12: Kết nối ống tiêu (peehole) liên kết trực tiếp ô trạng thái với các cổng

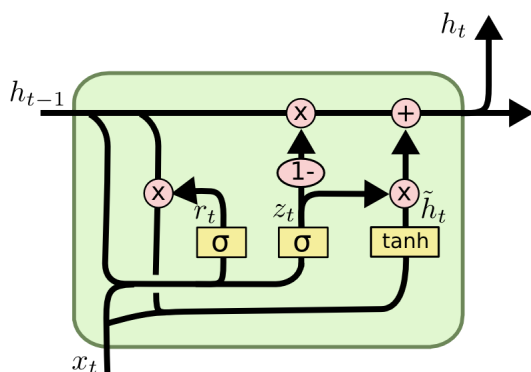
Một biến thể khác là sử dụng cặp đôi cổng vào và cổng ra. Thay vì quyết định riêng rẽ bỏ qua thông tin nào và thêm mới thông tin nào, chúng ta sẽ quyết định chúng đồng thời. Các thông tin chỉ bị quên khi chúng ta muốn cập nhập vào một vài thông tin mới.



$$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$$

**Hình 13: Cấu trúc điều chỉnh thêm mới và bỏ qua thông tin đồng thời**

Một dạng biến thể khá mạnh khác của LSTM là cổng truy hồi đơn vị (*Gated Recurrent Unit - GRU*) (<https://arxiv.org/pdf/1406.1078v3.pdf>) được giới thiệu bởi Cho, et al. (2014). Nó kết hợp cổng quên và cổng vào thành một cổng đơn gọi là cập nhật (*update gate*). Nó cũng nhập các ô trạng thái và trạng thái ẩn và thực hiện một số thay đổi khác. Kết quả của mô hình đơn giản hơn nhiều so với mô hình LSTM chuẩn, và đã trở nên khá phổ biến.



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

**Hình 14: Cấu trúc cổng truy hồi đơn vị (GRU - Gated Recurrent Unit)**

Chỉ có một số lượng nhỏ những biến thể LSTM là đáng chú ý. Rất nhiều những biến thể khác, như kiến trúc cổng sâu RNN (<https://arxiv.org/pdf/1508.03790v2.pdf>) (*Depth Gated RNN*) của Yao, et al. (2015) hay kiến trúc đồng hồ RNN (<https://arxiv.org/pdf/1402.3511v1.pdf>) (*Clockword RNN*) của Koutnik, et al. (2014) nhằm giải quyết vấn đề phụ thuộc dài hạn (*long - term dependencies*).

Vậy những biến thể nào là tốt nhất? Greff, et al. (2015) thực hiện một so sánh biến thể LSTM (<https://arxiv.org/pdf/1503.04069.pdf>) và nhận thấy rằng tất cả chúng đều giống nhau. Jozefowicz, et al. (2015) đã thử nghiệm hơn mười nghìn kiến trúc RNN, tìm thấy một số hoạt động tốt hơn LSTM trên một số nhiệm vụ nhất định. Trong khi đó Jozefowicz, et al. (2015) thực hiện kiểm tra trên hơn 1000 kiến trúc RNN (<http://proceedings.mlr.press/v37/jozefowicz15.pdf>) khác nhau và nhận thấy một số hoạt động tốt hơn so với LSTM trong một vài tác vụ cụ thể.

## 7. Kết luận

Trước đó, tôi đã đề cập đến những kết quả đáng chú ý mà mọi người đang đạt được với RNNs. Về cơ bản tất cả những điều này đều đạt được bằng cách sử dụng LSTM. Chúng thực sự làm việc tốt hơn rất nhiều cho hầu hết các nhiệm vụ!

Được viết dưới dạng một hệ phương trình, LSTM trông khá là phức tạp và có phần hằn lăm. Thông qua các bước diễn giải tuần tự nguyên lý hoạt động của nó tôi hi vọng sẽ khiến chúng trở nên dễ tiếp cận hơn.

LSTM là một bước đột phá lớn mà ở đó chúng ta đã khắc phục được những hạn chế ở RNN đó là khả năng phụ thuộc dài hạn. Một số kĩ thuật học Attention gần đây được kết hợp với LSTM đã tạo ra những kết quả khá bất ngờ trong các tác vụ dịch máy cũng như phân loại nội dung, trích lọc thông tin,.... Các mô hình dịch máy của google đã ứng dụng kiểu kết hợp này trong các bài toán dịch thuật của mình và đã cải thiện được nội dung bản dịch một cách đáng kể.

Top

## 8. Thực hành mô hình sinh từ tự động

### 8.1 Xây dựng mô hình trên level kí tự

Sau đây ta sẽ áp dụng mô hình LSTM trong việc dự báo từ tiếp theo của một đoạn hoặc câu văn dựa vào bối cảnh của từ là những từ liền trước nó.

Dữ liệu được sử dụng là bộ truyện alice ở xứ sở kỳ diệu

([https://gist.github.com/phillipj/4944029/raw/75ba2243dd5ec2875f629bf5d79f6c1e4b5a8b46/alice\\_in\\_wonderland.tx](https://gist.github.com/phillipj/4944029/raw/75ba2243dd5ec2875f629bf5d79f6c1e4b5a8b46/alice_in_wonderland.tx)) đã được nhà xuất bản publish nên không vi phạm bản quyền. Mô hình dự báo sẽ được xây dựng trên level kí tự.

Bên dưới, chúng ta sẽ đọc dữ liệu và chuyển các kí tự về in thường để giảm thiểu kích thước bộ mã hóa mà vẫn đảm bảo được nội dung văn bản. Dữ liệu được lưu trong kernel là file `wonderland.txt`.

```
1 import numpy
2 from keras.models import Sequential
3 from keras.layers import Dense, Dropout, LSTM
4 from keras.callbacks import ModelCheckpoint
5 from keras.utils import np_utils
6 import os
7
8 filename = '../input/wonderland.txt'
9 raw_text = open(filename).read().lower()
```

Một dictionary gồm 59 kí tự được sử dụng để mã hóa các kí tự trong bộ truyện. Key của các kí tự là số thứ tự của chúng trong dictionary. Trong sơ đồ thiết kế mạng nơ ron một kí tự được mã hóa bằng một vector đơn vị sao cho phần từ 1 sẽ xuất hiện tại vị trí của key trong bộ từ điển và 0 là các phần tử còn lại.

```
1 chars = sorted(list(set(raw_text)))
2 char_to_int = dict((c, i) for i, c in enumerate(chars))
3 print('number of letters: ', len(char_to_int))
4 print(char_to_int)

1 number of letters: 59
2 {'\n': 0, ' ': 1, '!': 2, '"': 3, '#': 4, '$': 5, '%': 6, '&': 7, '(': 8, ')': 9, '*'
```

Chúng ta nhận thấy rằng mục đích chỉ là dự báo từ tiếp theo do đó cần lọc bỏ những kí tự không quyết định đến nghĩa của 1 từ chẳng hạn như các dấu đặc biệt #, \$, \*, @, / . Như vậy, sẽ cần một bước chuẩn hóa dữ liệu nhằm giảm thiểu nhiễu và số lượng các khả năng có thể ở đầu ra. Điều này sẽ giúp cải thiện chất lượng và độ chính xác trong dự báo của mô hình đáng kể. Việc chuẩn hóa sẽ bao gồm như sau:

1. Chỉ giữ lại các kí tự chữ cái vì chúng có ảnh hưởng đến nội dung của 1 từ.
2. Chỉ giữ lại các dấu câu là ., !, ? vì chúng thể hiện các loại câu khác nhau và sẽ ảnh hưởng đến từ tiếp theo khi dự báo. Chẳng hạn nếu dấu câu là ? thì khả năng cao từ tiếp theo sẽ là yes hoặc no . Dấu câu là . thì từ tiếp theo có thể là một đại từ nhân xưng i, you, we, they, he, she, it .
3. Giữ lại các dấu , ' ' vì chúng giúp tách các từ và tách các thành phần câu.
4. Chuẩn hóa lại các chữ số về 1 chữ số duy nhất là 0 vì các con số là ngẫu nhiên và không dự báo được. Chúng ta chỉ có thể dự báo ở vị trí nào có khả năng là số.
5. Các kí tự nằm ngoài số liệt kê trên sẽ đưa vào nhóm unk tức unknown.

```
1 import string
2 string.ascii_lowercase
3 # string.digits
4 # string.punctuation
5 chars_new = list(string.ascii_lowercase) + ['0', '.', ',', '!', '?', 'unk']
6 chars_to_int = dict((v, k) for k, v in enumerate(chars_new))
7 int_to_chars = dict((k, v) for k, v in enumerate(chars_new))
8 print('character to int:', chars_to_int)
9 print('int to character:', int_to_chars)
10 # def _clean_char(text):
11 #     return 1
```

Top

```

1 character to int: {'a': 0, 'b': 1, 'c': 2, 'd': 3, 'e': 4, 'f': 5, 'g': 6, 'h': 7, 'i': 8}
2 int to character: {0: 'a', 1: 'b', 2: 'c', 3: 'd', 4: 'e', 5: 'f', 6: 'g', 7: 'h', 8: 'i'}

```

```

1 n_chars = len(raw_text)
2 n_vocab = len(chars_new)
3 print('Total characters: ', n_chars)
4 print('Total Vocab: ', n_vocab)

```

```

1 Total characters: 163693
2 Total Vocab: 33

```

Như vậy sau chuẩn hóa văn bản của chúng ta sẽ bao gồm 163693 từ và 33 kí tự. Tiếp theo là một hàm chuyển hóa một câu thành một vector chỉ số các kí tự.

```

1 def _encode_sen(text):
2     text = text.lower()
3     sen_vec = []
4     for let in text:
5         if let in chars_new[:-1]:
6             idx = chars_to_int[let]
7         else:
8             idx = chars_to_int['unk']
9         sen_vec.append(idx)
10    return sen_vec
11
12 x_test = _encode_sen('Alice is a wonderful story. #')
13 print(x_test)

```

```

1 [0, 11, 8, 2, 4, 29, 8, 18, 29, 0, 29, 22, 14, 13, 3, 4, 17, 5, 20, 11, 29, 18, 19, 1]

```

```

1 def _decode_sen(vec):
2     text = []
3     for i in vec:
4         let = int_to_chars[i]
5         text.append(let)
6     text = ''.join(text)
7     return text
8
9 _decode_sen(x_test)

```

```

1 'alice is a wonderful story. unk'

```

Để đồng nhất độ dài đầu vào cho mô hình cần tạo ra các chuỗi kí tự (window input) với độ dài là 100. Mục đích của chúng ta là dự báo kí tự tiếp theo từ 100 kí tự đầu vào. Mỗi một phiên dự báo window input sẽ được tịnh tiến lên 1 kí tự để thu được các kí tự dự báo liên tiếp nhau và từ đó ghép lại thành một câu hoàn chỉnh.

```

1 # prepare the dataset of input to output pairs encoded as integers
2 seq_length = 100
3 dataX = []
4 dataY = []
5 for i in range(0, n_chars - seq_length, 1):
6     # Lấy ra 100 kí tự liền trước
7     seq_in = raw_text[i:i + seq_length]
8     # Lấy ra kí tự liền sau 100 kí tự đó
9     seq_out = raw_text[i + seq_length]
10    dataX.append(_encode_sen(seq_in))
11    dataY.append(_encode_sen(seq_out)[0])
12    n_patterns = len(dataX)
13    print("Total Patterns: ", n_patterns)

```

Top



1 Total Patterns: 163593

Để có thể đưa vào mô hình LSTM, đầu vào  $x$  cần được chuẩn hóa thành một ma trận 3 chiều `samples, time steps, features`. Trong đó:

1. `samples`: Số lượng mẫu đầu vào (tức số lượng cửa sổ window 100 length).
2. `time steps`: Độ dài của cửa sổ window chính là số lượng các vòng lặp khi được trải phẳng ở hình 2 cấu trúc trải phẳng mạng nơ ron. Trong mô hình này `time steps` = 100.
3. `features`: Số lượng chiều được mã hóa của đầu vào. Trong mô hình LSTM, mỗi một từ hoặc kí tự (tùy theo chúng ta làm việc với level nào) thường được mã hóa theo 2 cách thông thường sau đây:
  - mã hóa theo one-hot encoding để một kí tự (ở bài thực hành này là kí tự) được biểu diễn bởi một véc tơ one-hot.
  - mã hóa theo giá trị véc tơ được lấy từ mô hình word embedding pretrain trước đó. Có thể là word2vec, fasttext, ELMo, BERT,... Số lượng các chiều theo level kí tự thường ít hơn so với level từ. Trong bài này để đơn giản ta không sử dụng một lớp embedding ở đầu để nhúng từ thành véc tơ mà sử dụng trực tiếp giá trị đầu vào là index của kí tự. Do đó số `features` = 1.

```
1 # reshape X to be [samples, time steps, features]
2 X_train = numpy.reshape(dataX, (n_patterns, seq_length, 1))
3 # normalize
4 X_train = X_train / float(n_vocab)
5 # one hot encode the output variable
6 y_train = np_utils.to_categorical(dataY)
7 print('X [samples, time steps, features] shape: ', X_train.shape)
8 print('Y shape: ', y_train.shape)
```

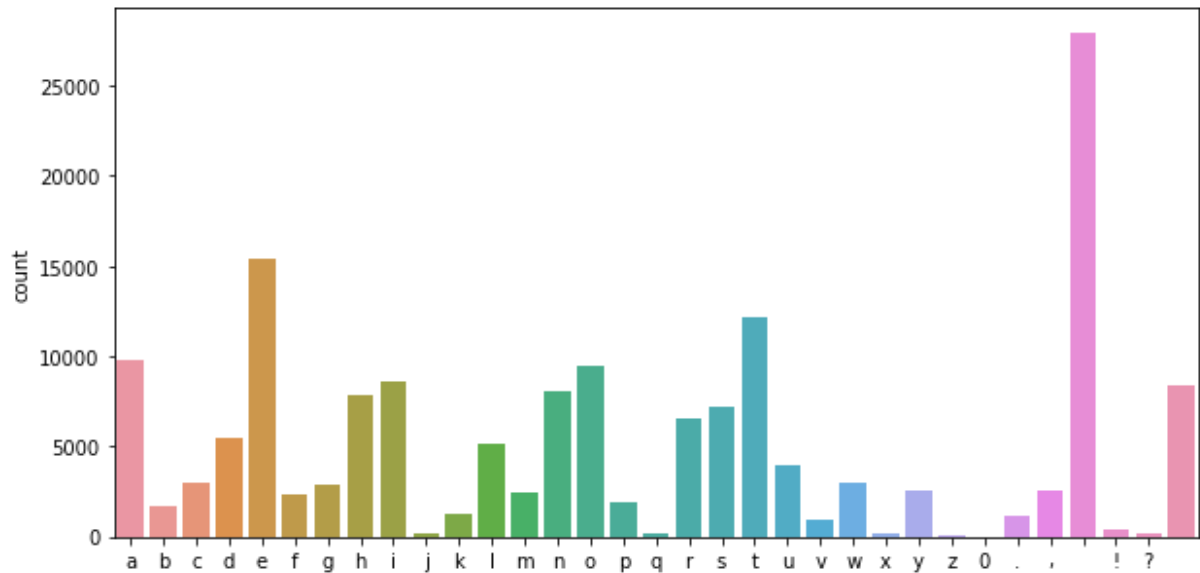
```
1 X [samples, time steps, features] shape: (163593, 100, 1)
2 Y shape: (163593, 33)
```

```
1 print(type(X_train))
2 print(type(y_train))
```

```
1 <class 'numpy.ndarray'>
2 <class 'numpy.ndarray'>
```

Thống kê số lượng các kí tự theo nhóm.

```
1 import seaborn as sn
2 import numpy as np
3 import matplotlib.pyplot as plt
4 plt.figure(figsize = (10, 5))
5 sn.countplot(np.array(dataY))
6 plt.xticks(np.arange(32), np.array(chars_new))
```



## 8.2. Mô hình LSTM

Xây dựng một kiến trúc model gồm một layer LSTM kết nối tới 1 layer Dropout và kết nối tới Dense layer ở cuối. Với mục đích là để chúng ta hiểu về lý thuyết LSTM nên tôi chọn mô hình có cấu trúc đơn giản nhất.

```

1 model = Sequential()
2 model.add(LSTM(256, input_shape = (X_train.shape[1], X_train.shape[2])))
3 model.add(Dropout(0.2))
4 model.add(Dense(y.shape[1], activation = 'softmax'))
5 model.compile(loss = 'categorical_crossentropy', optimizer = 'adam', metrics = ['accu
6 model.summary()

```

```

1
2 Layer (type)                Output Shape                Param #
3 =====
4 lstm_5 (LSTM)                (None, 256)                 264192
5
6 dropout_5 (Dropout)          (None, 256)                 0
7
8 dense_5 (Dense)              (None, 33)                  8481
9 =====
10 Total params: 272,673
11 Trainable params: 272,673
12 Non-trainable params: 0
13

```

```

1 filepath = 'weights-improvement-{epoch:02d}-{loss:.4f}.hdf5'
2 checkpoint = ModelCheckpoint(filepath, monitor = 'val_acc', verbose = 1, save_best_on
3 callback_list = [checkpoint]

```

```

1 model.fit(X_train, y_train, epochs = 5, batch_size = 128, validation_split=0.33, call

```

```

1 Train on 109607 samples, validate on 53986 samples
2 Epoch 1/5
3 109607/109607 [=====] - 408s 4ms/step - loss: 2.8991 - acc:
4
5 Epoch 00001: val_acc improved from -inf to 0.18686, saving model to weights-impro
6 Epoch 2/5
7 109607/109607 [=====] - 406s 4ms/step - loss: 2.7600 - acc:
8
9 Epoch 00002: val_acc improved from 0.18686 to 0.21489, saving model to weights-impro
10 Epoch 3/5
11 109607/109607 [=====] - 404s 4ms/step - loss: 2.6864 - acc:
12
13 Epoch 00003: val_acc improved from 0.21489 to 0.21520, saving model to weights-impro
14 Epoch 4/5
15 109607/109607 [=====] - 405s 4ms/step - loss: 2.6287 - acc:
16
17 Epoch 00004: val_acc improved from 0.21520 to 0.23499, saving model to weights-impro
18 Epoch 5/5
19 109607/109607 [=====] - 403s 4ms/step - loss: 2.5696 - acc:
20
21 Epoch 00005: val_acc improved from 0.23499 to 0.24291, saving model to weights-impro
22
23
24
25
26
27 <keras.callbacks.History at 0x7fa09d77fd30>

```

Dự báo kết quả từ tiếp theo từ một tập hợp kí tự đầu vào.

```

1 import numpy as np
2 base_word = 'Alice was beginning to get very tired of sitting by her sister on the b
3
4 def _predict_let(text, len_sen = 1):
5     text_for = []
6     for i in range(len_sen):
7         x_input = np.array(_encode_sen(text)[-100:])/float(n_vocab)
8         if x_input.shape[0] < 100:
9             x_input = np.concatenate((np.zeros(100-x_input.shape[0]), x_input), axis
10             x_input = np.expand_dims(np.expand_dims(x_input, -1), 0)
11             y_prob = model.predict(x_input)
12             y_let = int_to_chars[np.argmax(y_prob, axis = 1)[0]]
13             text = text + y_let
14     return text[len_sen:]
15
16 _predict_let(base_word, 100)

```

```

1 'nd the and the and the and the and the and the and the and the and t'

```

## 8.3. Kiến trúc BiLSTM.

Bên dưới là kiến trúc mạng LSTM 2 chiều có khả năng đọc đầu vào theo chiều từ trái qua phải và từ phải qua trái.

```

1  from keras.models import Sequential, Model
2  from keras.layers import Dense, Activation, Dropout
3  from keras.layers import LSTM, Input, Bidirectional
4  from keras.optimizers import Adam
5  from keras.callbacks import EarlyStopping
6  from keras.metrics import categorical_accuracy
7
8  #import spacy, and spacy english model
9  # spacy is used to work on text
10 import spacy
11 nlp = spacy.load('en')
12
13 #import other libraries
14 import numpy as np
15 import random
16 import sys
17 import os
18 import time
19 import codecs
20 import collections
21 from six.moves import cPickle
22
23 #define parameters used in the tutorial
24 data_dir = '../input' # data directory containing raw texts
25 # save_dir = 'save' # directory to store trained NN models
26 file_list = os.listdir('../input')
27 vocab_file = os.path.join("words_vocab.pkl")
28 sequences_step = 1 #step to create sequences

```

```

1  def create_wordlist(doc):
2      wl = []
3      for word in doc:
4          if word.text not in ("\n", "\n\n", '\u2009', '\xa0'):
5              wl.append(word.text.lower())
6      return wl
7
8  wordlist = []
9
10 for file_name in file_list:
11     input_file = os.path.join(data_dir, file_name)
12     #read data
13     with codecs.open(input_file, "r") as f:
14         data = f.read()
15
16     #create sentences
17     doc = nlp(data)
18     wl = create_wordlist(doc)
19     wordlist = wordlist + wl

```

Top

```

1  batch_size = 32 # minibatch size
2  num_epochs = 50 # number of epochs
3
4  callbacks=[EarlyStopping(patience=4, monitor='val_loss'),
5              ModelCheckpoint(filepath=save_dir + "/" + 'my_model_gen_sentences.{epoch:
6                  monitor='val_loss', verbose=0, mode='auto', period=2)]
7  #fit the model
8  history = md.fit(X, y,
9                  batch_size=batch_size,
10                 shuffle=True,
11                 epochs=num_epochs,
12                 callbacks=callbacks,
13                 validation_split=0.1)
14
15  #save the model
16  md.save(save_dir + "/" + 'my_model_generate_sentences.h5')batch_size = 32 # minibatch
17  num_epochs = 50 # number of epochs
18
19  callbacks=[EarlyStopping(patience=4, monitor='val_loss'),
20              ModelCheckpoint(filepath=save_dir + "/" + 'my_model_gen_sentences.{epoch:
21                  monitor='val_loss', verbose=0, mode='auto', period=2)]
22  #fit the model
23  history = md.fit(X, y,
24                  batch_size=batch_size,
25                  shuffle=True,
26                  epochs=num_epochs,
27                  callbacks=callbacks,
28                  validation_split=0.1)
29
30  #save the model
31  md.save(save_dir + "/" + 'my_model_generate_sentences.h5')batch_size = 32 # minibatch
32  num_epochs = 50 # number of epochs
33
34  callbacks=[EarlyStopping(patience=4, monitor='val_loss'),
35              ModelCheckpoint(filepath=save_dir + "/" + 'my_model_gen_sentences.{epoch:
36                  monitor='val_loss', verbose=0, mode='auto', period=2)]
37  #fit the model
38  history = md.fit(X, y,
39                  batch_size=batch_size,
40                  shuffle=True,
41                  epochs=num_epochs,
42                  callbacks=callbacks,
43                  validation_split=0.1)
44
45  #save the model
46  md.save(save_dir + "/" + 'my_model_generate_sentences.h5')len(wordlist[3])

```

```

1  # count the number of words
2  word_counts = collections.Counter(wordlist)
3
4  # Mapping from index to word : that's the vocabulary
5  vocabulary_inv = [x[0] for x in word_counts.most_common()]
6  vocabulary_inv = list(sorted(vocabulary_inv))
7
8  # Mapping from word to index
9  vocab = {x: i for i, x in enumerate(vocabulary_inv)}
10 words = [x[0] for x in word_counts.most_common()]
11
12 #size of the vocabulary
13 vocab_size = len(words)
14 print("vocab size: ", vocab_size)
15
16 #save the words and vocabulary
17 with open(os.path.join(vocab_file), 'wb') as f:
18     cPickle.dump((words, vocab, vocabulary_inv), f)

```

Top

```

1  #create sequences
2  sequences = []
3  next_words = []
4  for i in range(0, len(wordlist) - seq_length, sequences_step):
5      sequences.append(wordlist[i: i + seq_length])
6      next_words.append(wordlist[i + seq_length])
7
8  print('nb sequences:', len(sequences))

1  X = np.zeros((len(sequences), seq_length, vocab_size), dtype=np.bool)
2  y = np.zeros((len(sequences), vocab_size), dtype=np.bool)
3  for i, sentence in enumerate(sequences):
4      for t, word in enumerate(sentence):
5          X[i, t, vocab[word]] = 1
6          y[i, vocab[next_words[i]]] = 1

1  print('X shape: ', X.shape)
2  print('y shape: ', y.shape)

1  def bidirectional_lstm_model(seq_length, vocab_size):
2      print('Build LSTM model.')
3      model = Sequential()
4      model.add(Bidirectional(LSTM(rnn_size, activation="relu"), input_shape=(seq_length, vocab_size)))
5      model.add(Dropout(0.6))
6      model.add(Dense(vocab_size))
7      model.add(Activation('softmax'))
8
9      optimizer = Adam(lr=learning_rate)
10     callbacks=[EarlyStopping(patience=2, monitor='val_loss')]
11     model.compile(loss='categorical_crossentropy', optimizer=optimizer, metrics=[categorical_accuracy])
12     print("model built!")
13     return model
14
15     rnn_size = 256 # size of RNN
16     seq_length = 100 # sequence length
17     learning_rate = 0.001 #learning rate
18
19     md = bidirectional_lstm_model(seq_length, vocab_size)
20     md.summary()

1  batch_size = 32 # minibatch size
2  num_epochs = 50 # number of epochs
3
4  callbacks=[EarlyStopping(patience=4, monitor='val_loss'),
5             ModelCheckpoint(filepath="/" + 'my_model_gen_sentences.{epoch:02d}-{val_loss:0.2f}.h5',
6                             monitor='val_loss', verbose=0, mode='auto', period=2)]
7  #fit the model
8  history = md.fit(X, y,
9                  batch_size=batch_size,
10                 shuffle=True,
11                 epochs=num_epochs,
12                 callbacks=callbacks,
13                 validation_split=0.1)
14
15  #save the model
16  md.save("/") + 'my_model_generate_sentences.h5')

```

## 9. Tài liệu

Top

1. Understanding - LSTMs - Christopher Olah (<http://colah.github.io/posts/2015-08-Understanding-LSTMs>)

2. Learning Long - Term Dependencies with Gradient Descent is Difficult (<http://ai.dinfo.unifi.it/paolo/ps/tnn-94-gradient.pdf>)
3. (*Gated Recurrent Unit - GRU*) (<https://arxiv.org/pdf/1406.1078v3.pdf>)
4. so sánh biến thể LSTM (<https://arxiv.org/pdf/1503.04069.pdf>)
5. kiến trúc cổng sâu RNN (*Depth Gated RNN*) (<https://arxiv.org/pdf/1508.03790v2.pdf>)