



Deep Learning

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir



<https://www.aparat.com/mehran.safayani>



https://github.com/safayani/deep_learning_course

Examples of sequence data

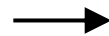
Speech recognition



“The quick brown fox jumped over the lazy dog.”

Music generation

∅



Sentiment classification

“There is nothing to like in this movie.”



DNA sequence analysis

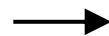
AGCCCCTGTGAGGAACTAG



AG**CCCCTGTGAGGAACT**AG

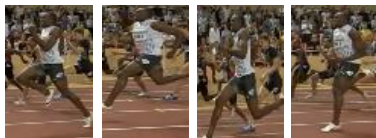
Machine translation

Voulez-vous chanter avec moi?



Do you want to sing with me?

Video activity recognition



Running

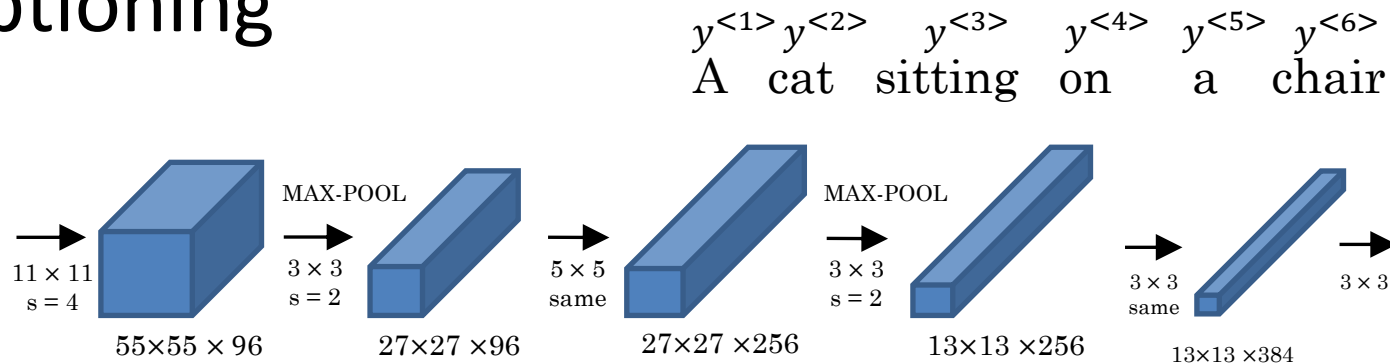
Name entity recognition

Yesterday, Harry Potter met Hermione Granger.

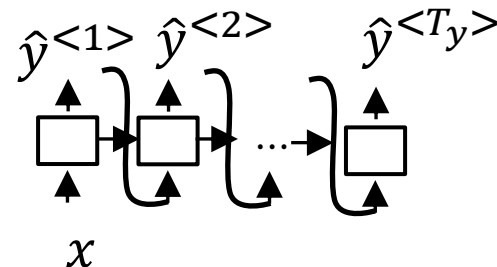
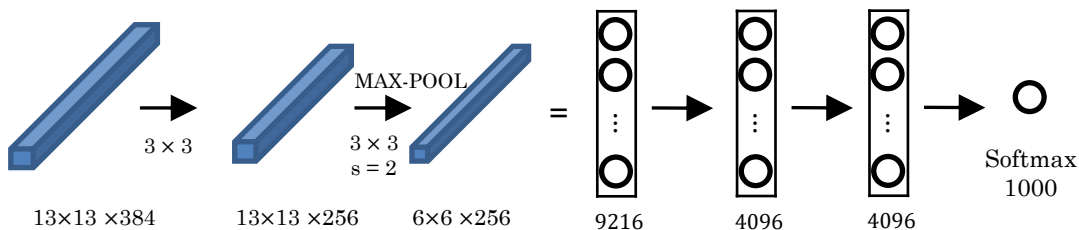


Yesterday, **Harry Potter** met **Hermione Granger**.

Image captioning



$y^{<1>}$ $y^{<2>}$ $y^{<3>}$ $y^{<4>}$ $y^{<5>}$ $y^{<6>}$
A cat sitting on a chair



[Mao et. al., 2014. Deep captioning with multimodal recurrent neural networks]

[Vinyals et. al., 2014. Show and tell: Neural image caption generator]

[Karpathy and Li, 2015. Deep visual-semantic alignments for generating image descriptions]

Recurrent Neural Networks

- Now, given this input X let's say that you want a model to operate Y that has one outputs per input word and the target output the design Y tells you for each of the input words is that part of a person's name.

- X: (Ali Ahmadi)and(Hassan Hamidi) invented a new drug.

$$X^{(1)} \quad X^{(2)} \quad X^{(3)} \quad \dots \quad \dots \quad X^{(t)} \quad \dots \quad \dots \quad X^{(9)} \quad T_x=9$$

- y: $\begin{matrix} 1 & 1 & 0 & 1 & & 1 & & 0 & 0 & 0 & 0 \\ y^{(1)} & y^{(2)} & y^{(3)} & & & & & & & & y^{(9)} \end{matrix} \quad T_y=9$

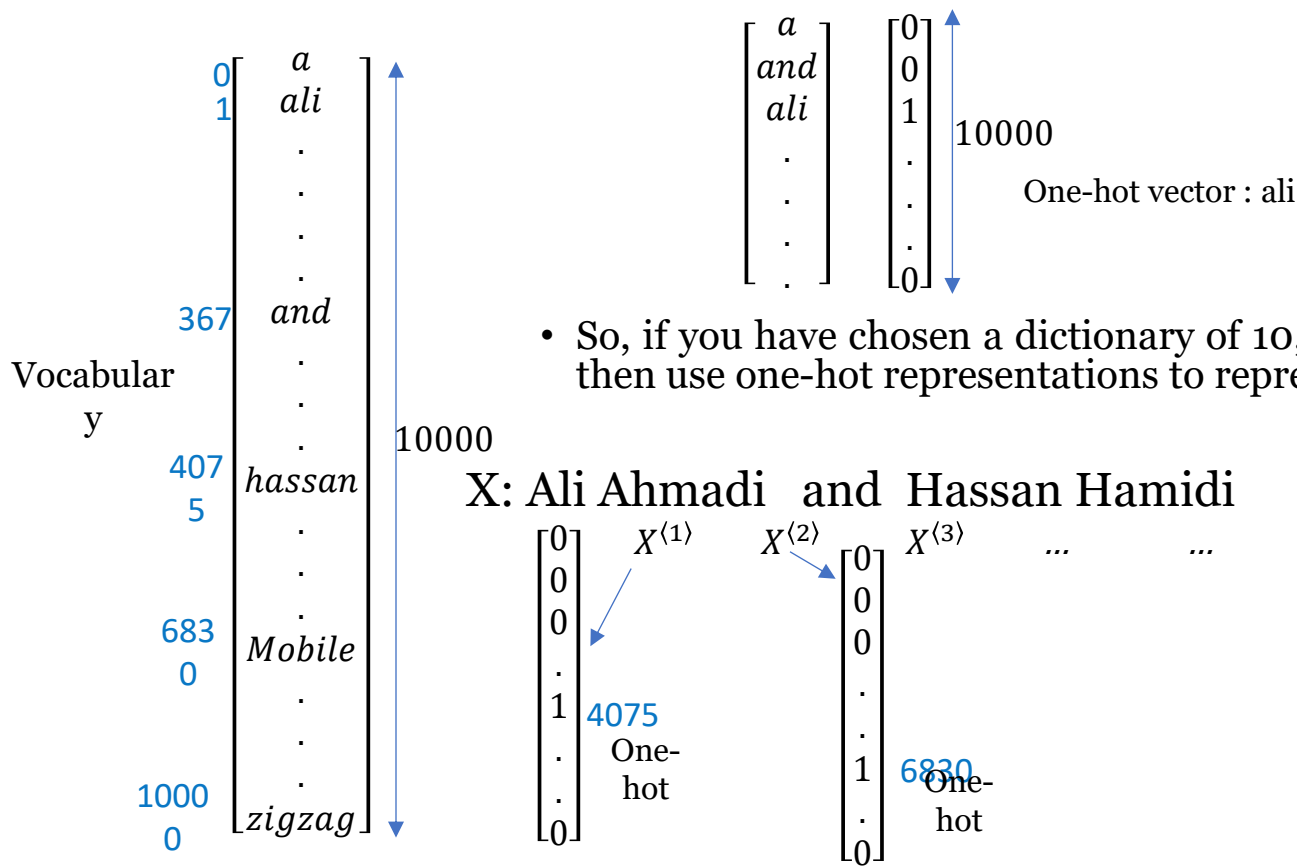
- $X^{(i)(t)}$

 داده آموزشی i ام
- $T_x^{(i)} = 9$
 $T_y^{(i)}$

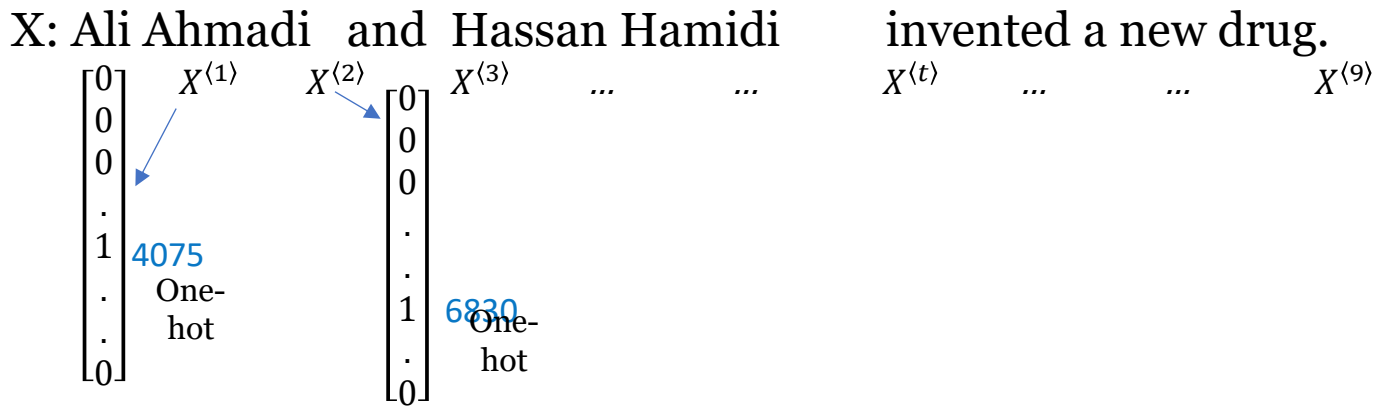
- This is our first serious foray into NLP or Natural Language Processing.

Representing words

- Dictionary: 30000, 50000



- So, if you have chosen a dictionary of 10,000 words, what you can do is then use one-hot representations to represent each of these words.

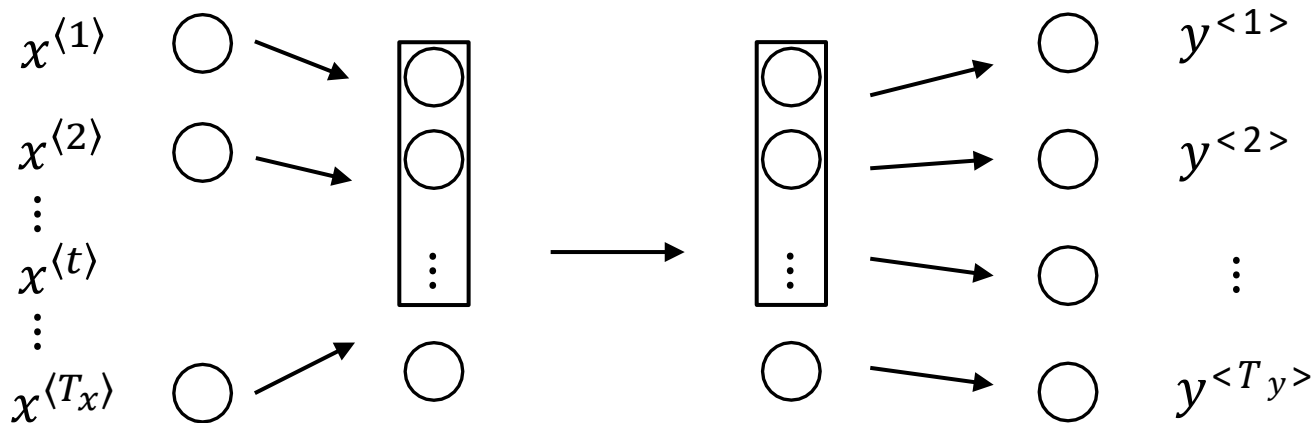


Representing words

- Utilizing a sequence model for supervised learning to map input X to output Y .
- Introduction of an "Unknown Word" token for handling out-of-vocabulary words.
- Describing a notation for training sets in sequence data.

Recurrent Neural Network Model

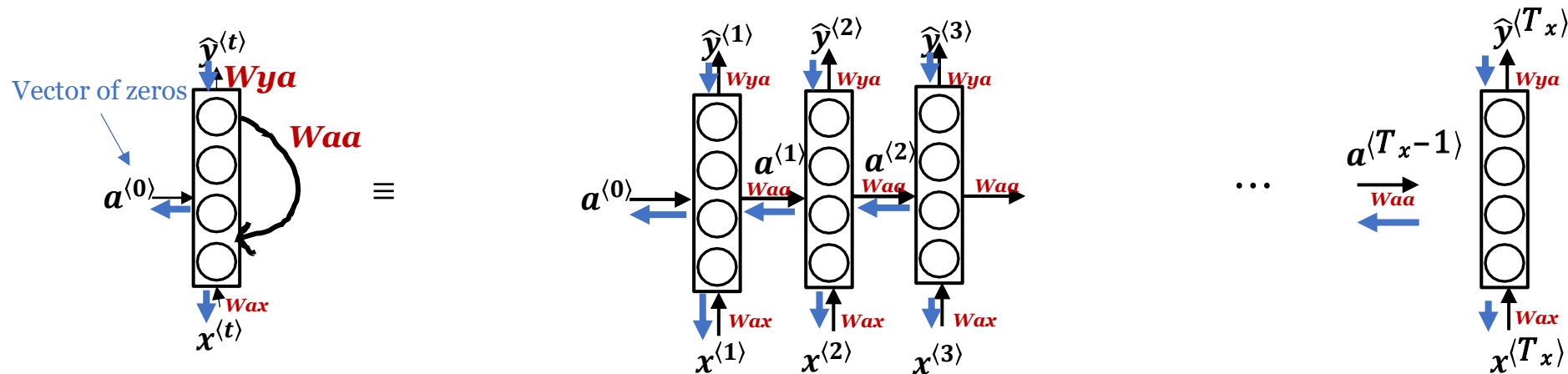
Why its not a standard network?



Problems:

- Inputs, outputs can be different lengths in different examples.
- Doesn't share features learned across different positions of text.

Recurrent Neural Networks



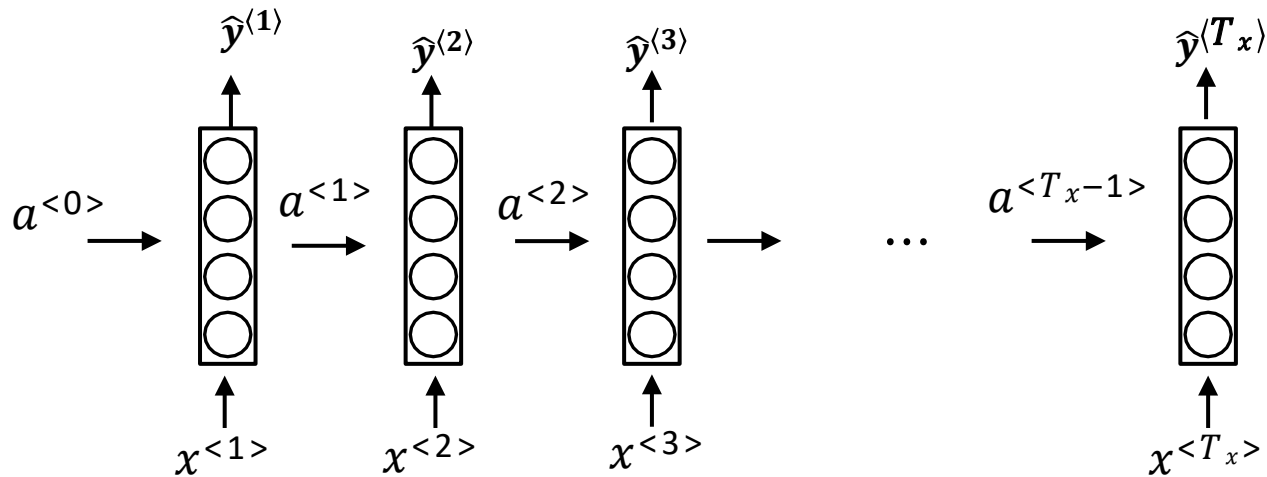
$$L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1 - y^{(t)}) \log(1 - \hat{y}^{(t)})$$

$$L(\hat{y}, y) = \sum_{t=1}^{T_x} L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

- Activation at time zero in neural networks is often initialized as a vector of zeros.
- Some researchers prefer initializing it as $a^{(0)}$ randomly.
- There are alternative methods to initialize the activation at time zero

Backward propagation through time

Forward Propagation



$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g_1(w_{aa}a^{<0>} + w_{ax}x^{<1>} + b_a) \leftarrow \text{tanh|Relu}$$

$$\hat{y}^{<1>} = g_2(w_{ya}a^{<1>} + b_y) \leftarrow \text{sigmoid}$$

$$a^{<t>} = g(w_{aa}a^{<t-1>} + w_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(w_{ya}a^{<t>} + b_y)$$

Simplified RNN notation

$$a^{(t)} = g(w_{aa}a^{(t-1)} + w_{ax}x^{(t)} + b_a)$$

Diagram annotations for the first equation:
 - w_{aa} has a blue arrow from $a^{(t-1)}$ to the product term, labeled $(100,100)$.
 - w_{ax} has a blue arrow from $x^{(t)}$ to the product term, labeled $(100,10000)$.
 - b_a has a blue arrow pointing to the sum, labeled 100 .
 - $a^{(t-1)}$ has a blue arrow pointing down to the value 100 .
 - $x^{(t)}$ has a blue arrow pointing down to the value 10000 .

$$a^{(t)} = g(w_a[a^{(t-1)}, x^{(t)}] + b_a)$$

$$\hat{y}^{(t)} = g(w_{ya}a^{(t)} + b_y)$$

$$\bullet a^{(t)} = g(w_a[a^{(t-1)'}', u^{(t)'}'] + b_a)$$

$$\begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix}$$

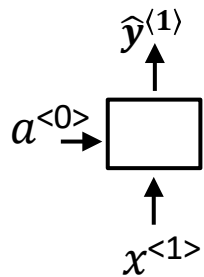
Diagram annotations for the vector equation:
 - $a^{(t-1)}$ has a blue double-headed vertical arrow next to it labeled 100 .
 - $x^{(t)}$ has a blue double-headed vertical arrow next to it labeled 10000 .
 - A blue double-headed vertical arrow to the right of the entire vector is labeled 10100 .

$$w_a = [w_{aa} : w_{ax}]$$

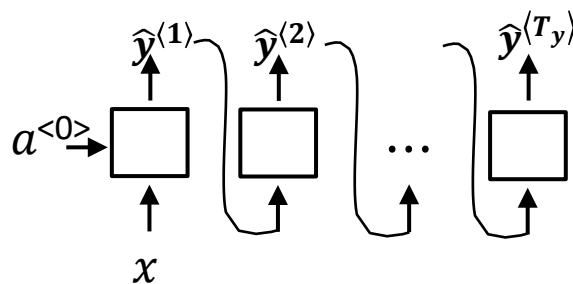
Diagram annotations for the weight matrix equation:
 - Above w_{aa} is a blue double-headed horizontal arrow labeled 100 .
 - Above w_{ax} is a blue double-headed horizontal arrow labeled 100 .
 - Between w_{aa} and w_{ax} is a blue double-headed horizontal arrow labeled 10000 .
 - Below the entire matrix is a blue double-headed horizontal arrow labeled $(100,10100)$.

$$[w_{aa} : w_{ax}] \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = w_{aa}a^{(t-1)} + w_{ax}x^{(t)}$$

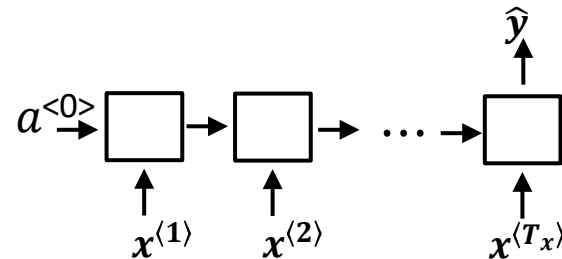
Summary of RNN types



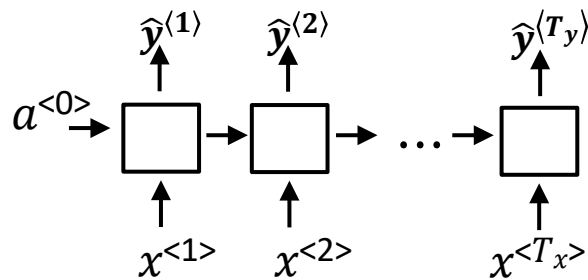
One to one



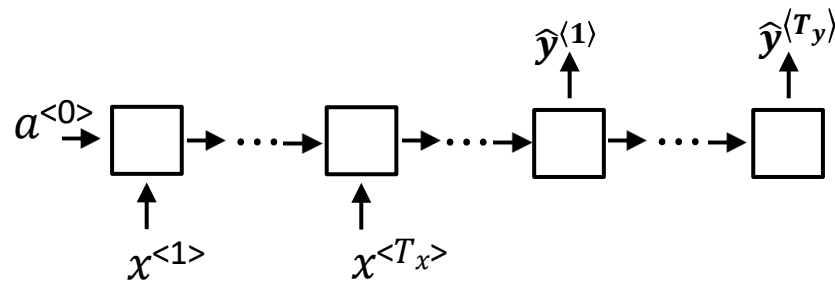
One to many



Many to one



Many to many



Many to many

What is language modelling?

Speech recognition

- The apple and pair salad.
- The apple and pear salad.
- $P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$
- $P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$
- $P(\text{sentence}) = P(y^{\langle 1 \rangle}, y^{\langle 2 \rangle}, \dots, y^{\langle T_y \rangle}) = \text{احتمال وقوع جمله}$

Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

Cats average 15 hours of sleep a day. \downarrow $\langle \text{EOS} \rangle$
 $y^{(1)}$ $y^{(2)}$ $y^{(3)}$... $y^{(9)}$

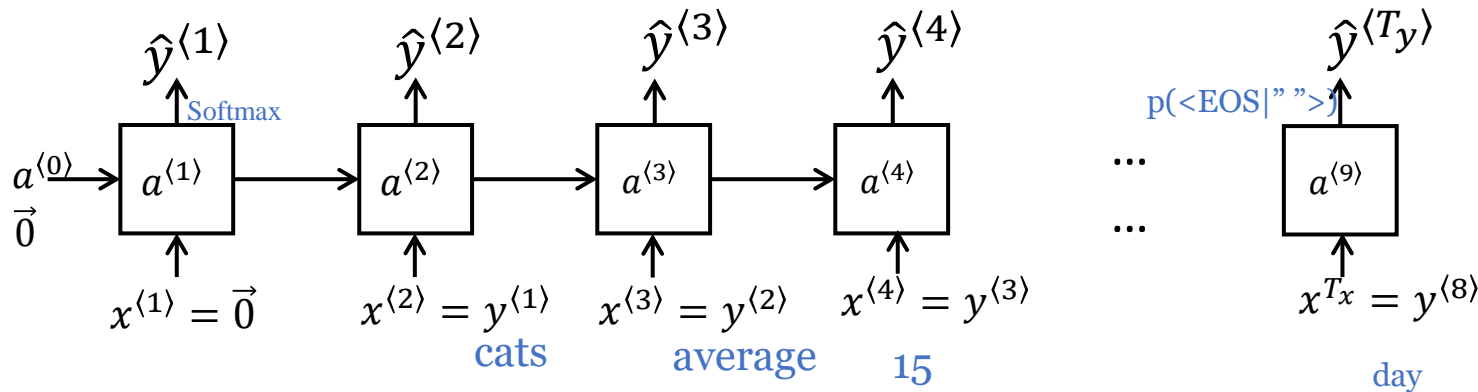
$$x^{(t)} = y^{(t-1)}$$

The Egyptian ~~Mau~~ is a breed of cat. $\langle \text{EOS} \rangle$
 $\langle \text{UNK} \rangle$

RNN model

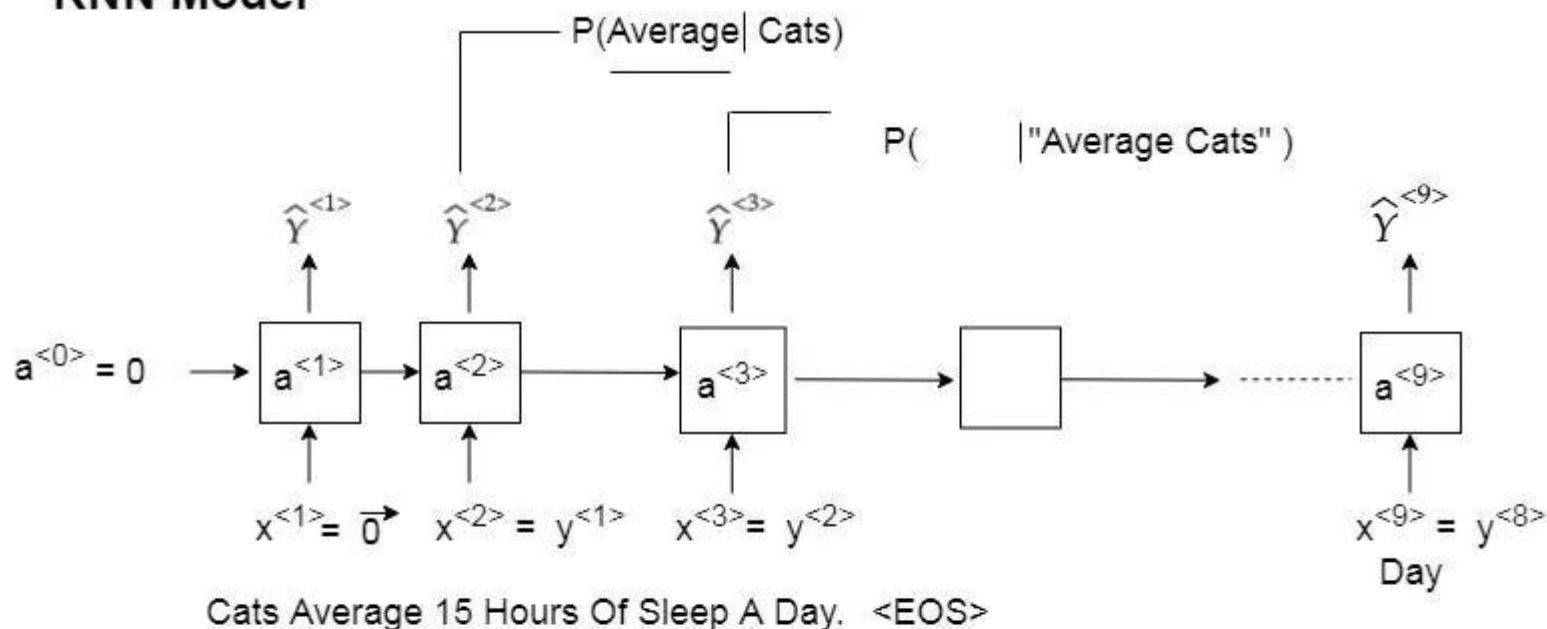
Cats average 15 hours of sleep a day.

- $p(a)$ $p(\text{cat})$ $p(\text{bread})$... $p(\text{Eat})$
- $p(a|\text{cats})$ $p(\text{cat}|\text{cats})$ $p(\text{bread}|\text{cats})$ $p(\text{average}|\text{cats})$ $p(\text{Eat}|\text{cats})$
- $P(a|\text{cats average})$ $p(\text{cats}|\text{cats average})$... $p(15|\text{cats average})$ $P(\text{eat}|\text{cats average})$



- $L(\hat{y}(t), y(t)) = -\sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$
- $L = \sum_t L^{(t)}(\hat{y}^{(t)}, y^{(t)})$

RNN Model

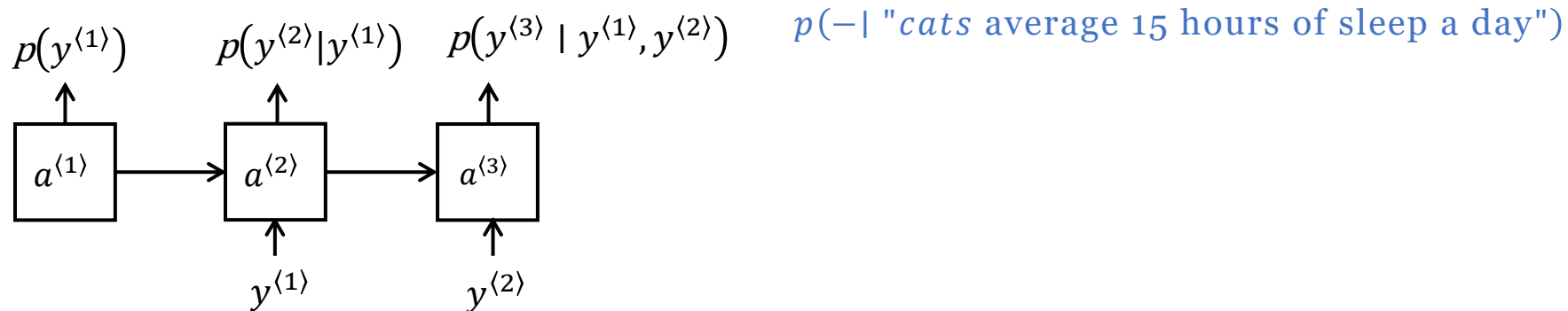


Cost Function $L(\hat{Y}^{<t>}, Y^{<t>}) = - \sum_i Y_i^{<t>} \log \hat{Y}_i^{<t>}$ ← SoftMax Loss Function

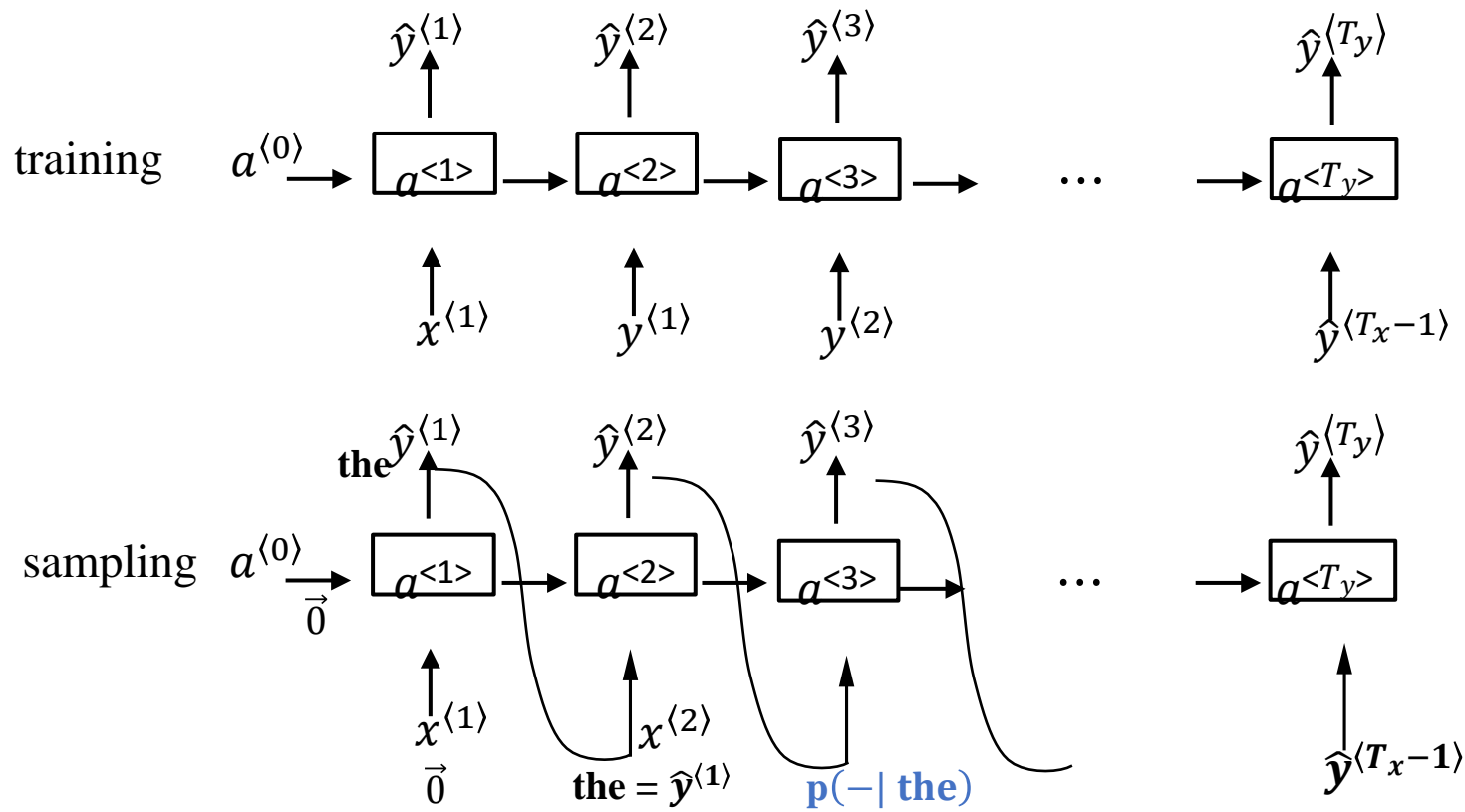
$L = \sum_t L^{<t>}(\hat{Y}^{<t>}, Y^{<t>})$ ← Overall Loss Function

RNN model

- Cats average 15 hours of sleep a day. <EOS>
- $p(y^{(1)}, y^{(2)}, y^{(3)}) = p(y^{(1)}) \cdot p(y^{(2)} | y^{(1)}) \cdot p(y^{(3)} | y^{(1)}, y^{(2)})$



Sampling a sequence from a trained RNN



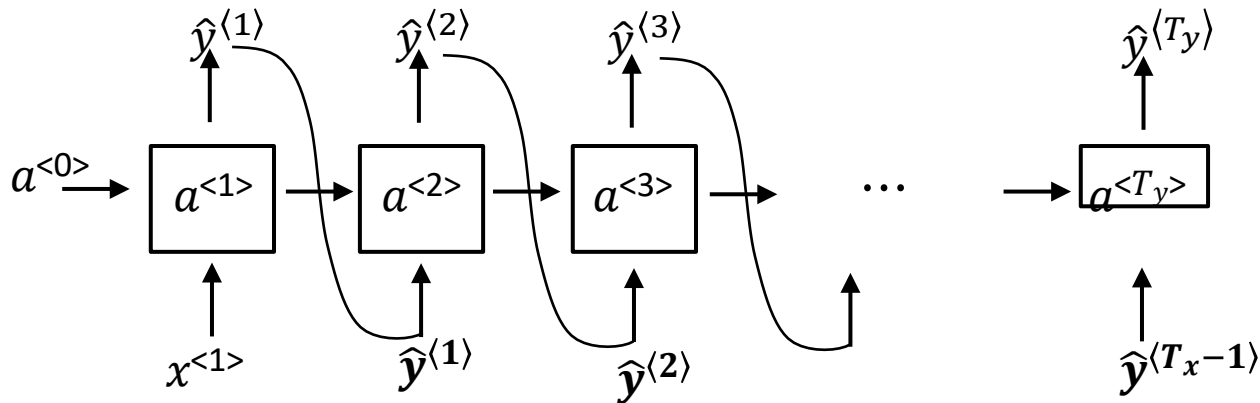
<EOS>

<UNK>

$p(a)$ $p(aaron)$... $p(zulu)$ $p(<UNK>)$ `np.random.choice`

Character-level language model

- Vocabulary = [a, aaron, ..., zulu, <UNK>]
- Vocabulary = [a, b, c, ... , z , \sqcup , . , ..., 0, ... , 9, A, ... ,Z]
- Cat average ... Mau
↑↑↑



Sequence generation

News

President Enrique Peña Nieto, announced
sench's sulk former coming football langston
paring.

"I was not at all surprised," said Hich Langston.

"Concussion epidemic", to be examined.

The gray football the told some and this has on
the UEFA icon, should money as.

Shakespeare

The mortal moon hath her eclipse in love.

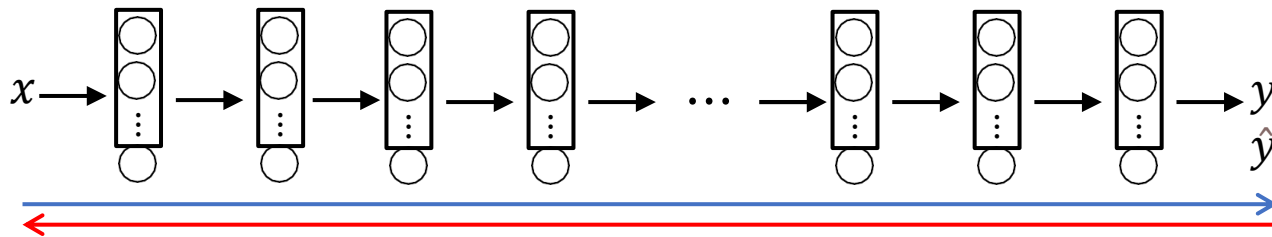
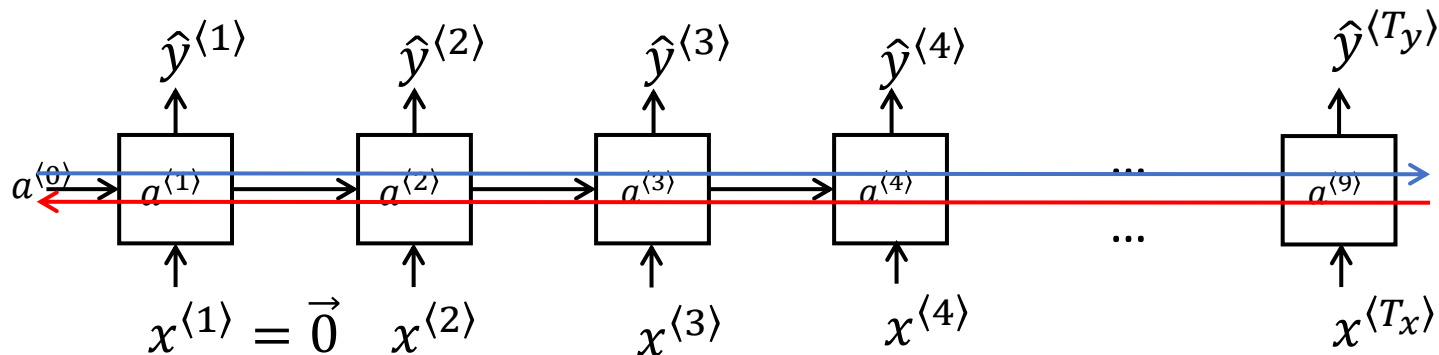
And subject of this thou art another this fold.

When better be my love to me see Sabl's.

For whose are ruse of mine eyes heaves.

Vanishing gradients with RNNs

- The cat which already ate bunch of food was full
- The cats ... were full



Exploding gradients.

NaN

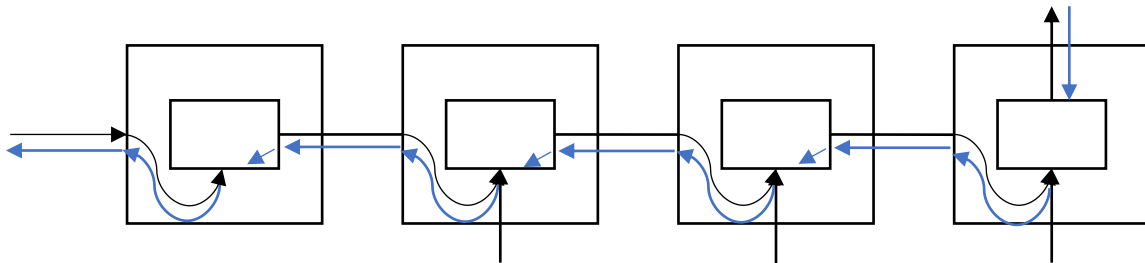
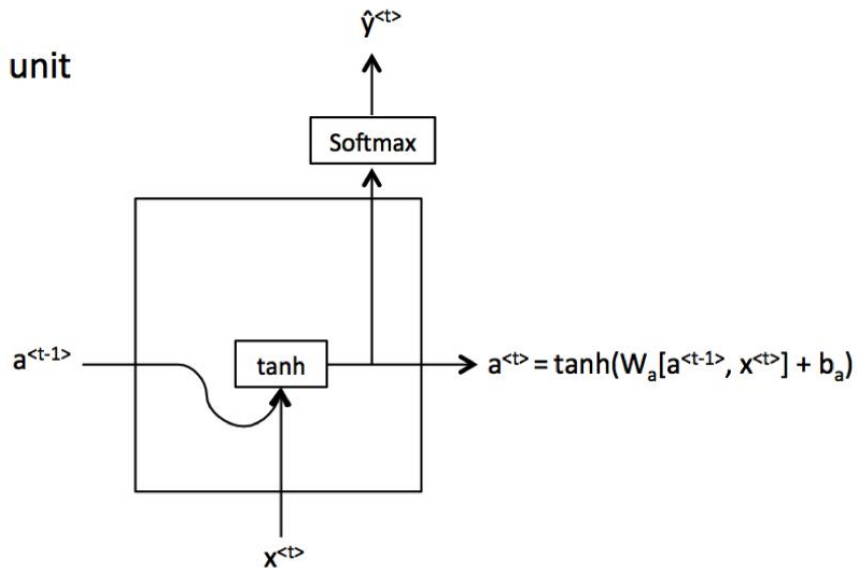
gradient clipping

RNN Unit

- $a^{(t)} = g(W_a[a^{(t-1)}, x^{(t)}] + b_a)$

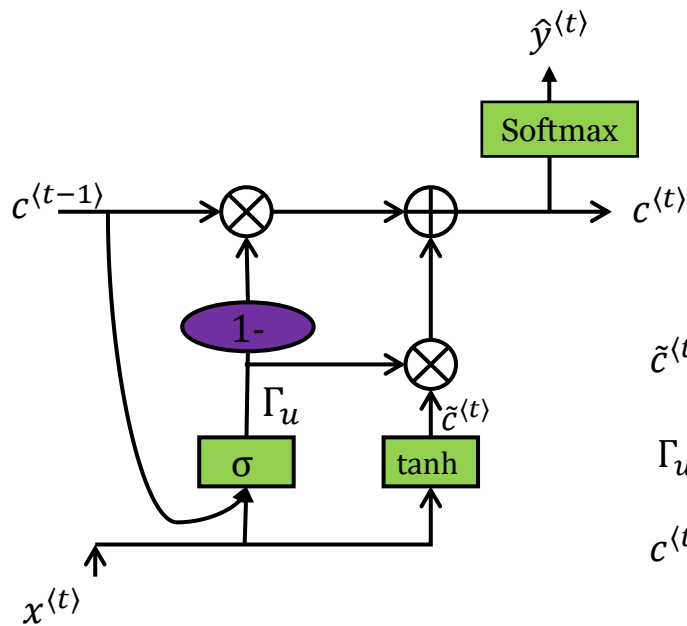
↑
tanh

RNN unit



GRU (simplified)

$\Gamma_u = 0$ $\Gamma_u = 1$ $\Gamma_u = 0$ $\Gamma_u = 0$ $\Gamma_u = 0$ \dots $\Gamma_u = 0$
 $c^{(t)} = 0$ $c^{(t)} = 1$ $c^{(t+1)} = 1$ $c^{(t+2)} = 1$ $c^{(t+3)} = 1$ \dots $c^{(t+n)} = 1$
 The **cat**, which already ate ..., **was** full.



C : memory cell

$$\tilde{c}^{(t)} = \tanh(w_c[c^{(t-1)}, x^{(t)}] + b_c)$$

$$\Gamma_u = \sigma(w_u[c^{(t-1)}, x^{(t)}] + b_u)$$

$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

Full GRU

- $\hat{c}^{(t)} = \tanh(w_c[\Gamma_r * c^{(t-1)}, x^{(t)}] + b_c)$
- $\Gamma_u = \sigma(w_u[c^{(t-1)}, x^{(t)}] + b_u)$
- $\Gamma_r = \sigma(w_r[c^{(t-1)}, x^{(t)}] + b_r)$
- $c^{(t)} = \Gamma_u * \hat{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$

