

Setting Up

$X^{(i)\leftarrow t}$ i^{th} training sample

$T_x^{(i)}$ Length of i^{th} training sample $0 \leq t \leq T_x^{(i)}$

$y^{(i)\leftarrow t}$

$T_y^{(i)}$

Vocabulary	One-hot Encoding
Word 1	0 0 1 0 0
Word 2	0 0 0 1 0
:	
Word 10,000	0 0 0 0 1

<unk>

Top 10,000 Most Common Words

Why not NN?

Standard NN Problems

- 1) Inputs, Outputs Diff lengths
- 2) Not Share features learned across diff Positions of text

Long-term Dependency

Activation Functions

$$a^{<0>} = \vec{0} \quad (\text{A Fake Vector As Start})$$

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

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

Simplified Vector Form

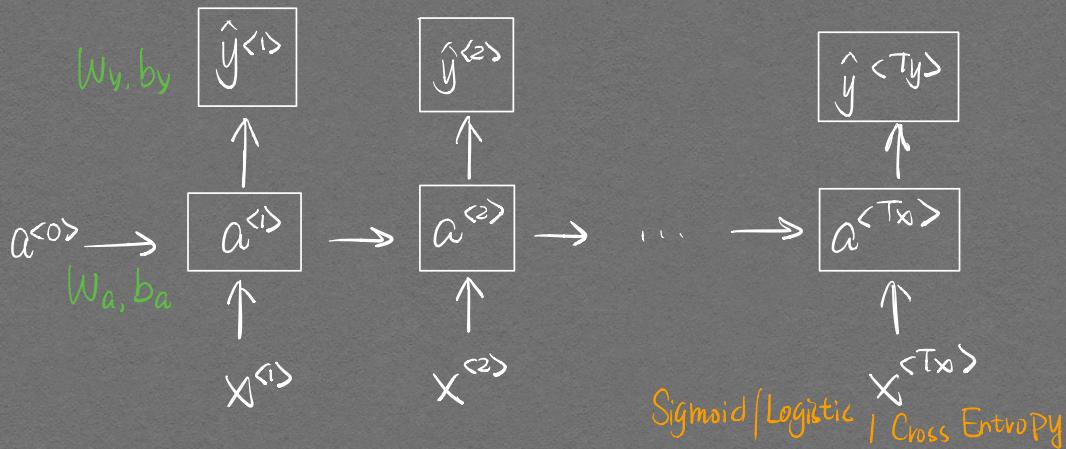
$$a^{<t>} = g(\underbrace{W_a [a^{<t-1>} X^{<t>}]}_{W_a : W_x} + b_a)$$

$$\begin{bmatrix} W_a & | & W_x \\ \hline \end{bmatrix} \begin{bmatrix} a^{<t-1>} \\ \hline \dots \\ \hline X^{<t>} \end{bmatrix} = W_a a^{<t-1>} + W_x X^{<t>}$$

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

Forward Propogation & Backward Propogation

$$\mathcal{L} = \sum (\mathcal{L}^{(1)}, \mathcal{L}^{(2)}, \dots, \mathcal{L}^{(T_y)})$$



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

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

Diff Types of RNNs

- Speech Recognition
- Music Generation
- Sentiment Classification
- DNA Sequence
- Machine Translation
- Video Activity Recognition
- Name Entity Recognition

$(T_x = T_y = 1)$ One To One

$T_x = T_y$ Many To Many

$T_x \neq T_y$ Encoder → Decoder

$T_x > T_y$
 $T_y = 1$ Many To One

$T_x < T_y$ One To Many
 $T_x = 1$

Language Model & Sequence Generation

$$P(\text{"The Apple & Pear Salad"})$$

<

$$P(\text{Sentence}) = ?$$

$$P(\text{"The Apple & Pear Salad"})$$

$$\text{Language Model } P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)})$$

Training Set: Large Corpus of English Text Sampling Novel Sequences

Tokenize

$$P(\text{word 1})$$

A Distribution Over Vocabulary

One-hot Encoding

$$P(\text{word 2})$$

np.random.choice

(Add <EOS> Tag)

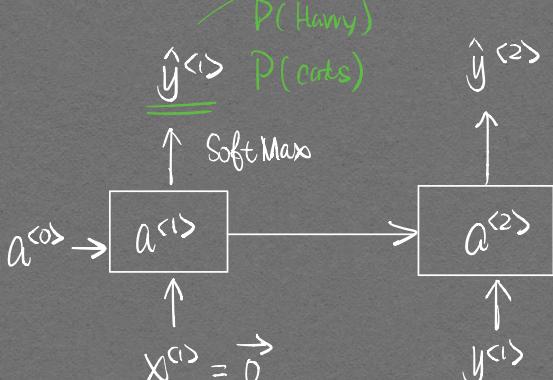
$$\vdots$$

while <UNK>:
Regenerate Word

<UNK> Tag / <OOV> Tag

$$P(\text{<UNK>})$$

$$P(\text{<EOS>} | \text{cate. ...})$$



$$\mathcal{L}(\hat{y}^{(t)}, y^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)} \quad \text{Soft Max Loss}$$

$$\mathcal{L} = \sum_t \mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

$$P(y_1, y_2, y_3) = P(y_1) \cdot P(y_2 | y_1) \cdot P(y_3 | y_1, y_2)$$

Language Model predicts Probabilities

Character Language Model

A B ... Z

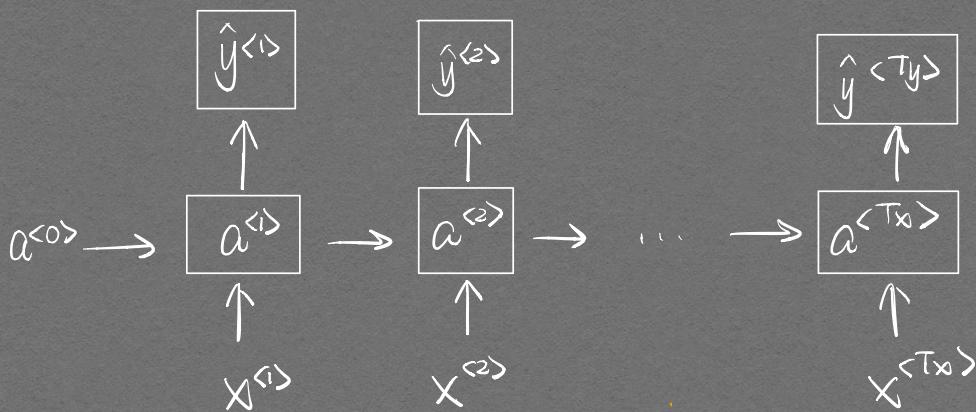
a b ... z

No <UNK> Word.

Computationally Expensive.

Specialized Applications

Vanishing Gradients With RNN



Back Propagation Doesn't work well on Earlier Inputs

Long-term Dependence → Long-term Memory

Requires Gradient Clipping ↗ Exploding Gradients

Gradients ↗ Increase Exponentially ↗ Blow-up Params → NaNs
↘ Decrease Exponentially ↗ Normalize Too Large Gradients ↗ Numerical Overflow

$\|g\|_1 \|g\|_2$

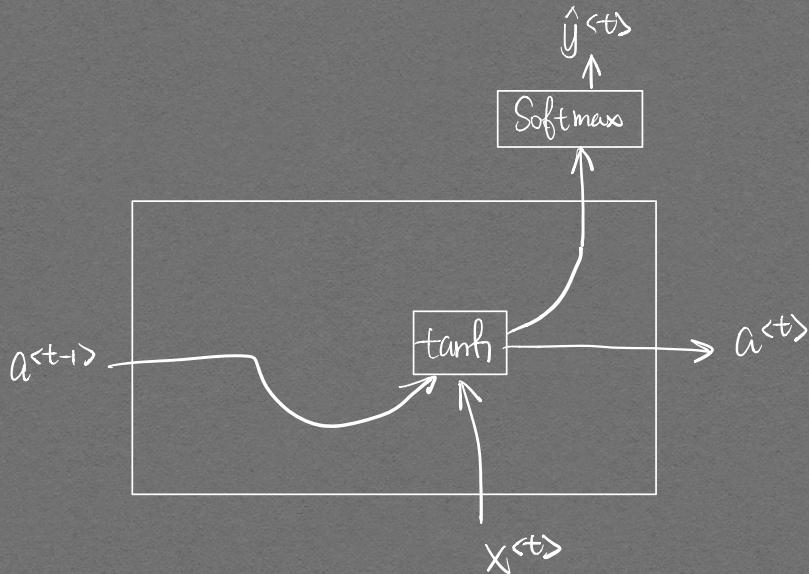
GRU ← Vanishing Gradients

LSTM

RNN Unit

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

$$\hat{y}^{t \rightarrow} = g(W_y a^{t \rightarrow} + b_y)$$



Simplified GRU Unit

C = memory cell

$$C^{t \rightarrow} = a^{t \rightarrow}$$

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

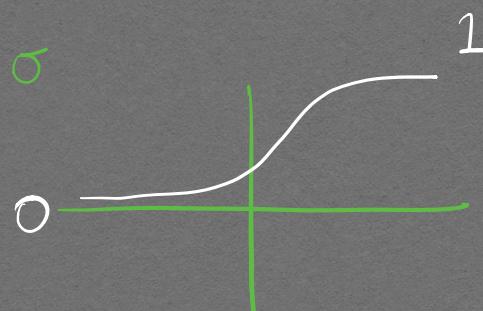
Candidate to Replace C

$$\Gamma_u = \sigma(W_u [c^{t-1}, x^{t \rightarrow}] + b_u) \quad \text{Easy To Set to 0}$$

(Gamma_u) $\stackrel{\text{Sigmoid 0/1}}{=} \text{"Update Gate"}$ Update or not

$$C^{t \rightarrow} = \Gamma_u \cdot \hat{C}^{t \rightarrow} +$$

$$(1 - \Gamma_u) \cdot \hat{C}^{t-1}$$



$$= \begin{cases} \hat{c}^{(t)}, P_u = 1 & \text{Update} \\ \hat{c}^{(t-1)}, P_u = 0 & \text{Don't Update} \end{cases}$$

or close to 0
very very small $1E-6$

✓ Maintaining Memory

Full GRU Unit

GRU

C = memory cell

$$C^{(t)} = a^{(t)}$$

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

Candidate to Replace C = relevance "how relevant $c^{(t-1)}$ is to compute $C^{(t)}$ "

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

(Gamma u) "Update Gate" Update or not A vector of dimension equal to

$$C^{(t)} = T_u \cdot \tilde{c}^{(t)} + \text{the number of hidden units}$$

$$(-T_u) \cdot \hat{c}^{(t-1)}$$

$$= \begin{cases} \hat{c}^{(t)}, P_u = 1 & \text{Update} \\ \hat{c}^{(t-1)}, P_u = 0 & \text{Don't Update} \end{cases}$$

$$T_r = \sigma(W_r [c^{(t-1)}, x^{(t)}] + b_r)$$

GRU

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

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

$$T_r = \sigma(W_r [c^{(t-1)}, x^{(t)}] + b_r)$$

LSTM

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

$$T_u = \sigma(W_u [\underline{a}^{(t-1)}, x^{(t)}] + b_u)$$

$$T_f = \sigma(W_f [\underline{a}^{(t-1)}, x^{(t)}] + b_f)$$



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

$$a^{(t)} = \underline{\tilde{c}^{(t)}}$$

Two Gates

Cheap

Scalable

forget

$$\Rightarrow P_o = \sigma (W_o [a^{(t-1)}, x^{(t)}] + b_o)$$

Output

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

$$a^{(t)} = \underline{\underline{\Gamma_o \cdot \tanh(\tilde{c}^{(t)})}}$$

Three Gates

Expensive

Flexible

Powerful

LSTM

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

Update

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

Forget

$$\Gamma_f = \sigma (W_f [a^{(t-1)}, x^{(t)}, c^{(t-1)}] + b_f)$$

Output

$$\Rightarrow P_o = \sigma (W_o [a^{(t-1)}, x^{(t)}, c^{(t-1)}] + b_o)$$

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

$$a^{(t)} = \underline{\underline{\Gamma_o \cdot \tilde{c}^{(t)}}}$$

"Peephole

Connection"

Bi-LSTM

Getting Information From Future
Acyclic Graph

W_y, b_y

$\hat{y}^{(1)}$

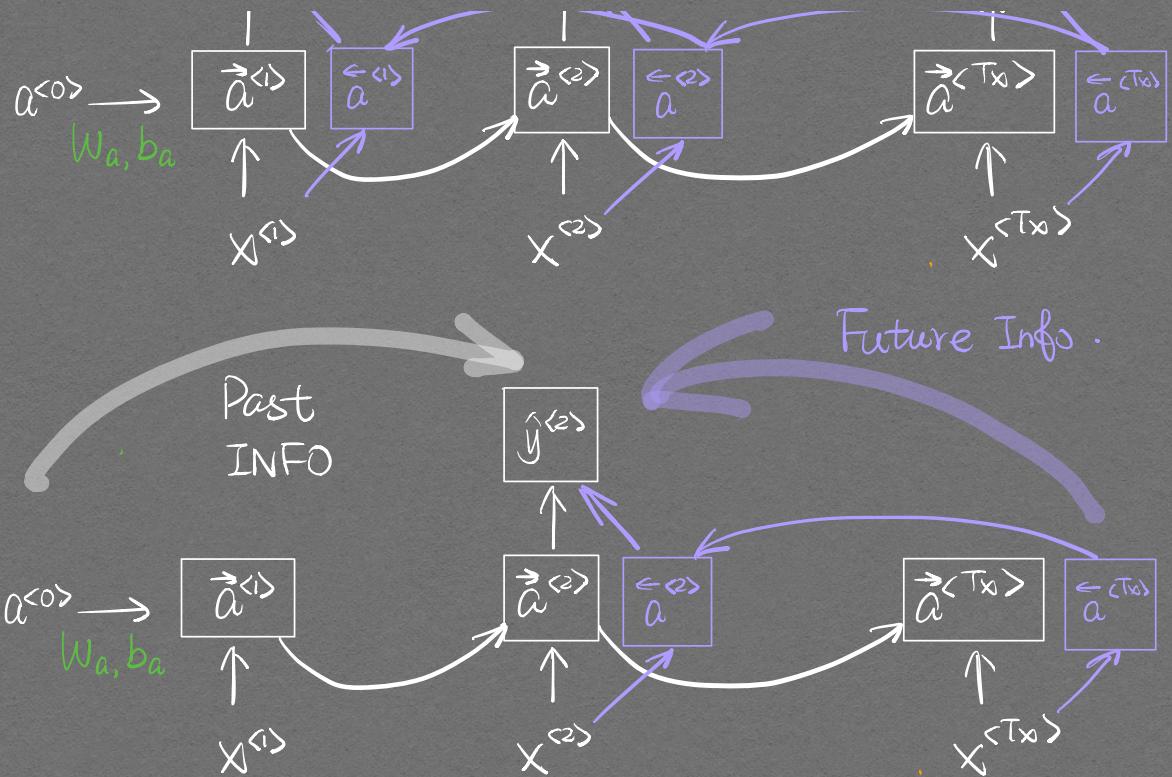
$\uparrow \nwarrow$

$\hat{y}^{(2)}$

$\uparrow \nwarrow$

$\hat{y}^{(T_y)}$

$\uparrow \nwarrow$



Deep-RNN

Standard NN Unroll Through Time

