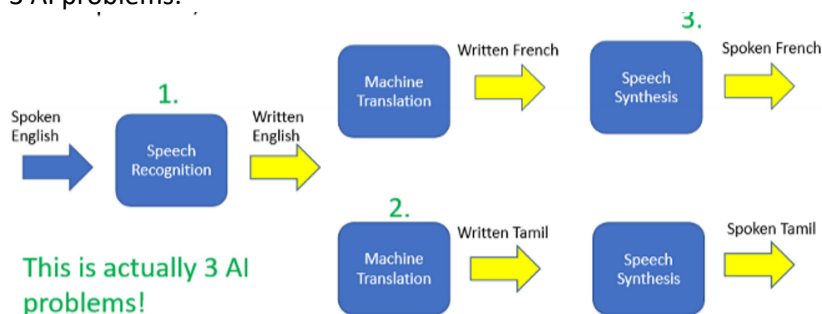# RNN and Natural Language Processing

Tuesday, December 14, 2021    6:52 PM
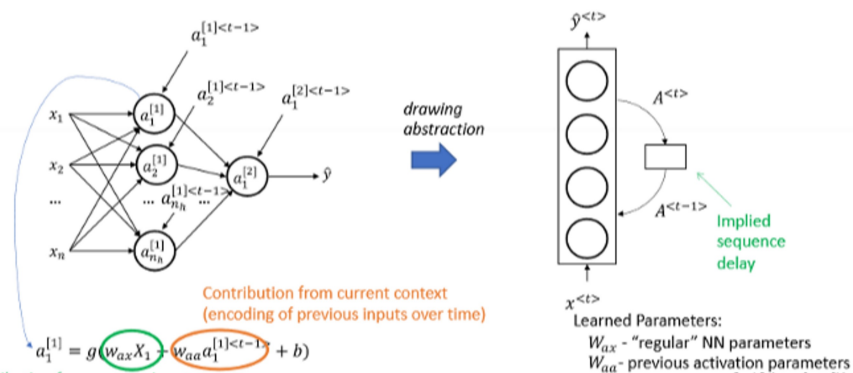
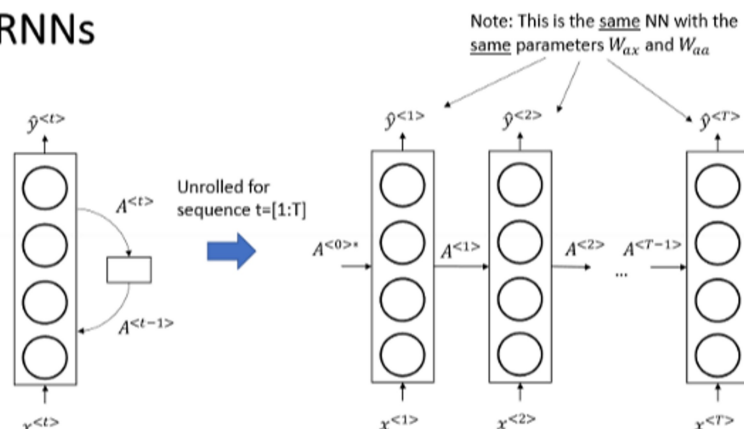**21. Recurrent Neural Networks**

**Translation**
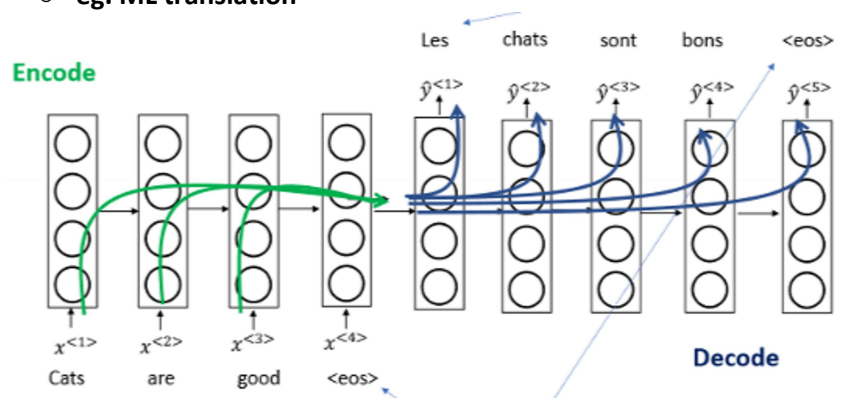- 3 AI problems:



- information in individual components of data and their ordering
- need to consider the context!
  - ○ eg. Michael Jordan made a basket
- **Giving ML Context**
  - ○ understand new info
  - ○ increasing inputs to system to reflect context (data from the past) doesn't scale!
  - ○ activations from previous step in sequence can be used to "bias" activations in next step
    - ▪ could simultaneously learn amount of context required while we learn input to output mappings



- **I/O Sequence Length Flexibility in RNNs**
  - ○ Flexibility around size of input/output data

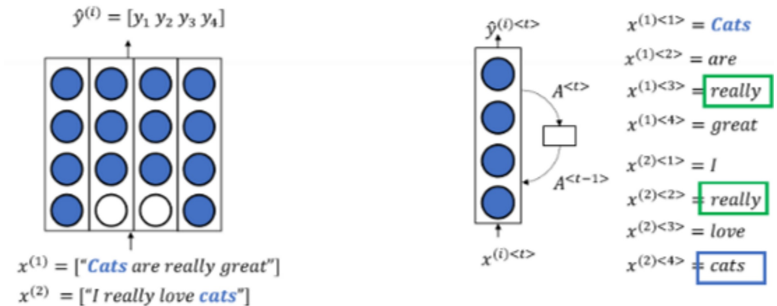| $\hat{y}^{<1>}$ | $\hat{y}^{<1>}$ $\hat{y}^{<2>}$ $\hat{y}^{<3>}$ | $\hat{y}^{<1>}$ | $\hat{y}^{<1>}$ $\hat{y}^{<2>}$ $\hat{y}^{<3>}$ |

$x^{<1>}$ | $x^{<1>}$ | $x^{<1>}$ $x^{<2>}$ $x^{<3>}$ | $x^{<1>}$ $x^{<2>}$ $x^{<3>}$

| a. one-to-one | b. one-to-many | *many-to-one* c. one-to-many | d. many-to-many |
| Image Classification | Image Captioning | Sentiment Classification | Machine Translation |
| | | | Brad Quinton, Scott Cl |

- ○ **eg. ML translation**



**Encode**

Les    chats    sont    bons    <eos>

$\hat{y}^{<1>}$    $\hat{y}^{<2>}$    $\hat{y}^{<3>}$    $\hat{y}^{<4>}$    $\hat{y}^{<5>}$

$x^{<1>}$    $x^{<2>}$    $x^{<3>}$    $x^{<4>}$

Cats    are    good    <eos>

**Decode**

- **RNN Feature Extraction**



$\hat{y}^{(i)} = [y_1\ y_2\ y_3\ y_4]$

$x^{(1)} = [\text{"Cats are really great"}]$
$x^{(2)} = [\text{"I really love cats"}]$

$\hat{y}^{(i)<t>}$

$A^{<t>}$

$A^{<t-1>}$

$x^{(i)<t>}$

$x^{(1)<1>} = Cats$
$x^{(1)<2>} = are$
$x^{(1)<3>} = really$
$x^{(1)<4>} = great$

$x^{(2)<1>} = I$
$x^{(2)<2>} = really$
$x^{(2)<3>} = love$
$x^{(2)<4>} = cats$

- ▪ eg. neuron learned to activate strongly to patterns that make an 11x11 swatch of facial hair
- ○ RNNs isolate elements of sequences like conv filters isolate regions of an image
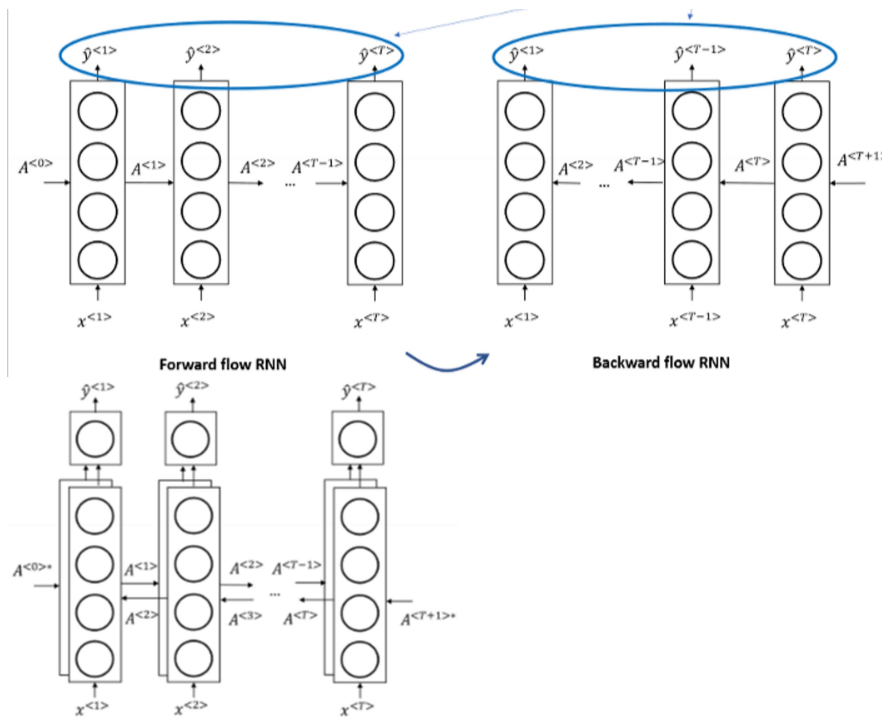- **Bidirectional RNN**
  - ○ can look forward and backward in time
  - ○ buffer the inputs long enough to consider context in 2 directions
    - ▪ context may only be defined later in the sentence



"The fair approach is to split the profits 50/50."

"The fair was in town for only one night."

The context is defined later in the sentence.    Brad Quir

  - ○ create forward flow RNN and backward flow RNN then combine outputs
  - ○ using BRNNs for each sentence = state of the art for NLP

**Forward flow RNN**    **Backward flow RNN**



### 23. Implementation of RNNs

### RNN Notation

- **Inputs:** $x^{(i)<t>}$ where $i = 1:m$ and $t = 1:T_x^{(i)}$

  Now, each input and output in the training example has a sequence length T.

- **Outputs:** $y^{(i)<t>}$ where $i = 1:m$ and $t = 1:T_y^{(i)}$

  Like before, we have $m$ training examples.

  - eg. for x: "cats are nice", want to make each word an element of our sequence
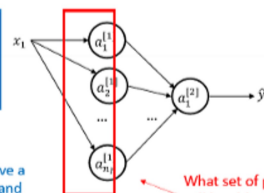
$$x^{<1>} = \text{"Cats"}, x^{<2>} = \text{"are"}, x^{<3>} = \text{"nice"}, x^{<4>} = \text{"."}$$

  - need to assign each word a number
  - can create an ordered dictionary, assign each word a num based on position in sequence
    - may add unintential bias, words are biased based on position in alphabet

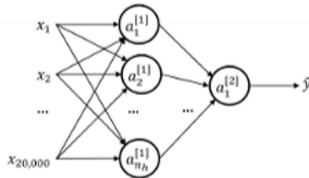| Word | Position |
|------|----------|
| a | 1 |
| aaron | 2 |
| ... | |
| cats | 520 |
| .... | |
| nice | 3024 |
| ... | |
| zulu | 19,999 |
| <eos>, "." | 20,000 |

"Cats": 520
"Are": 97
"Nice": 3024
"": 20,000

Why should "a" and "aaron" have a similar magnitude, but not "a" and "erin", or "tom"? Words are biased together based on their position in the alphabet!

What set of parameters will be able to create a distinct activation for each word?

  - need normalized representation and less compressed
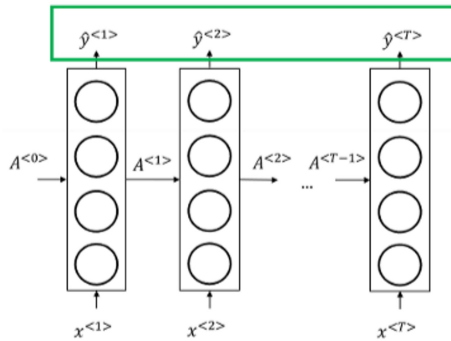    - instead use one-hot encoding, no order bias

- **Unknown words**
    - ○ one more vector element called unknown word ("UKW")
    - ○ works as long as all important words NP task uses are included in vocabulary

**RNN Loss Function**
- define loss function to optimize later
- define overall loss to be sum of each element of sequence
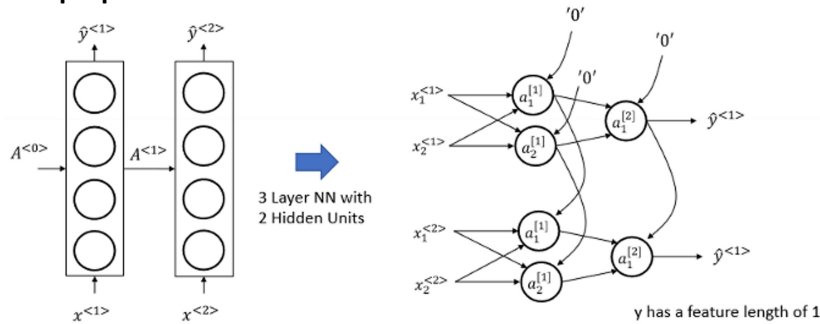- using one-hot encodings, re-use Cross Entropy Loss for each element



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

And, given we are using one-hot encodings, we can re-use Cross Entropy Loss for each element:

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

- **Backprop in RNNs**



Let's imagine we have a $T_y = T_x = 2$

x has a feature length of 2

y has a feature length of 1

3 Layer NN with 2 Hidden Units

1. calculate y_hat using computation graph
2. determine the loss
3. update each parameter (using partial derivative of cost)
4. repeat until J<target
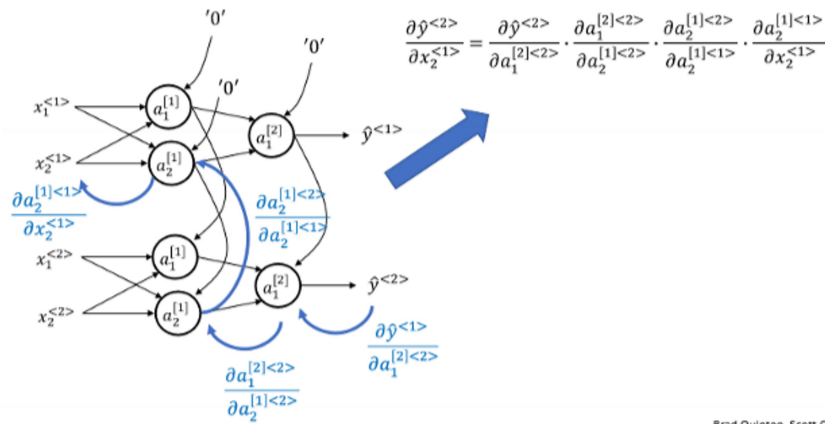


Step 1: Calculate $\hat{y}$ using computation graph.

Step 2: Determine the loss.

$$J = -\frac{1}{m}\left(\sum_{i=1}^{m} y^i\log(\hat{y}^{(i)}) + \sum_{i=1}^{m}(1 - y^{(i)})\log(1 - \hat{y}^{(i)})\right)$$

Step 3: Update *each* parameter (Using the partial derivative of cost).

$$w = w - \alpha\frac{\partial J}{\partial w}; b = b - \alpha\frac{\partial J}{\partial b}$$
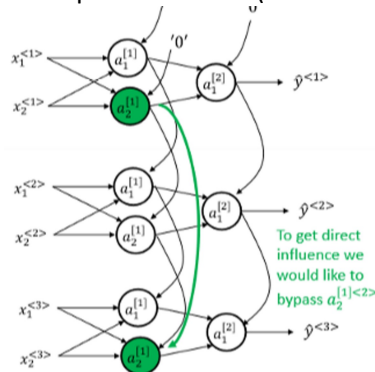
- **Partial Derivatives on RNN**

$$\frac{\partial \hat{y}^{<2>}}{\partial x_2^{<1>}} = \frac{\partial \hat{y}^{<2>}}{\partial a_1^{[2]<2>}} \cdot \frac{\partial a_1^{[2]<2>}}{\partial a_2^{[1]<2>}} \cdot \frac{\partial a_2^{[1]<2>}}{\partial a_2^{[1]<1>}} \cdot \frac{\partial a_2^{[1]<1>}}{\partial x_2^{<1>}}$$

- ○ **Vanishing gradients in RNNs**
    - ▪ as sequences get long, can be difficult to enable earlier elements to correctly influence later outputs
    - ▪ influence of early elements of the sequences are "washed out" by repeat multiply accumulates in each iteration
    - ▪ activation unit for RNN:

$$a_1^{[1]<t>} = g(w_{ax}X_1 + w_{aa}a_1^{[1]<t-1>} + b)$$

    - ▪ repeated applying same RNN activation units. To bypass current activation, we hold previous value (add memory to RNN)



- **Simplified Gated Recurrent Unit (GRU)**
    - ○ get a value between 0 and 1 based on learned param and standard RNN unit inputs

$$\Gamma_\mu = \sigma(w_{\mu x}X_1 + w_{\mu a}a_1^{[1]<t-1>} + b_\mu)$$

    - ○ Use function to decide whether keep previous activation or update

$$a_1^{[1]<t>} = \Gamma_\mu * \tilde{a}_1^{[1]<t>} + (1 - \Gamma_\mu) * a_1^{[1]<t-1>}$$

    - ○ standard activation calculation now a candidate:

$$\tilde{a}_1^{[1]<t>} = g(w_{ax}X_1 + w_{aa}a_1^{[1]<t-1>} + b)$$

        - ▪ similar to creating a mux function
- **Long Short Term Memory (LSTM)**
    - ○ maintian influence of value across many sequences
    - ○ manage vanishing gradient problem in RNN
    - ○ create new var c, expression of internal memory
    - ○ gating: update, forget, output

$$\Gamma_\mu = \sigma(w_{\mu x}X_1 + w_{\mu a}a_1^{[1]<t-1>} + b_\mu) \qquad \text{Update}$$

$$\Gamma_f = \sigma(w_{fx}X_1 + w_{fa}a_1^{[1]<t-1>} + b_f) \qquad \text{Forget}$$

$$\Gamma_o = \sigma(w_{ox}X_1 + w_{oa}a_1^{[1]<t-1>} + b_o) \qquad \text{Output}$$

- use funcs to manage internal memory and output vals
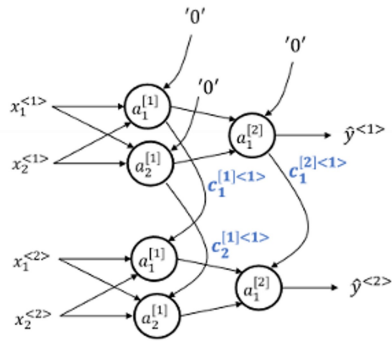- eg. create candidate memory:

$$\tilde{c}_1^{[1]<t>} = g(w_{ax}X_1 + w_{aa}a_1^{[1]<t-1>} + b)$$

- update internal memory, update and forget

$$c_1^{[1]<t>} = \Gamma_\mu * \tilde{c}_1^{[1]<t>} + \Gamma_f * c_1^{[1]<t-1>}$$

- create output separate from internal memory

$$a_1^{[1]<t>} = \Gamma_o * \tanh(c_1^{[1]<t>})$$



- structure of many sequential data sets, key element critical for period of time, then no longer relevant

## 24. Natural Language Processing (NLP)