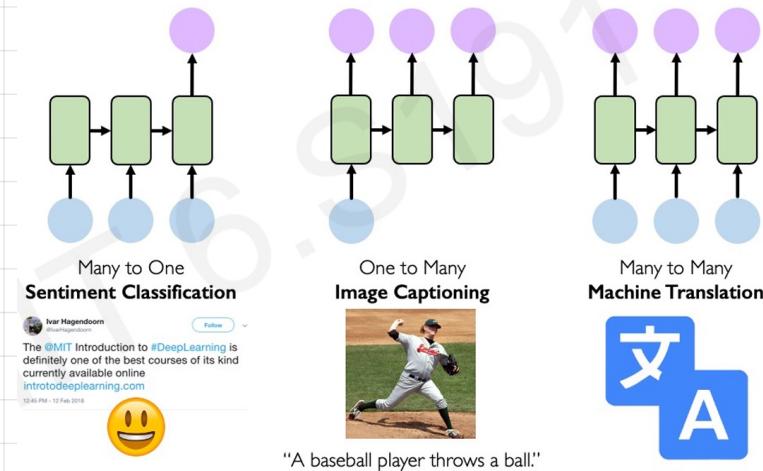


Sequential data

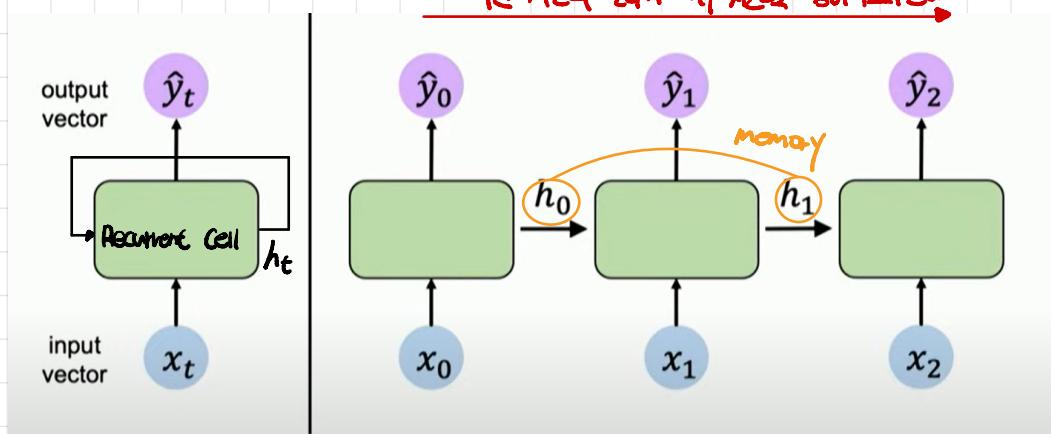
- 음파, 텍스트의 흐름, 퍼스트, 주가...

Sequential Modeling Application



Neurons with Recurrence

이전 시간의 출력과 이후 시간의 출력이 서로 영향을 미친다.

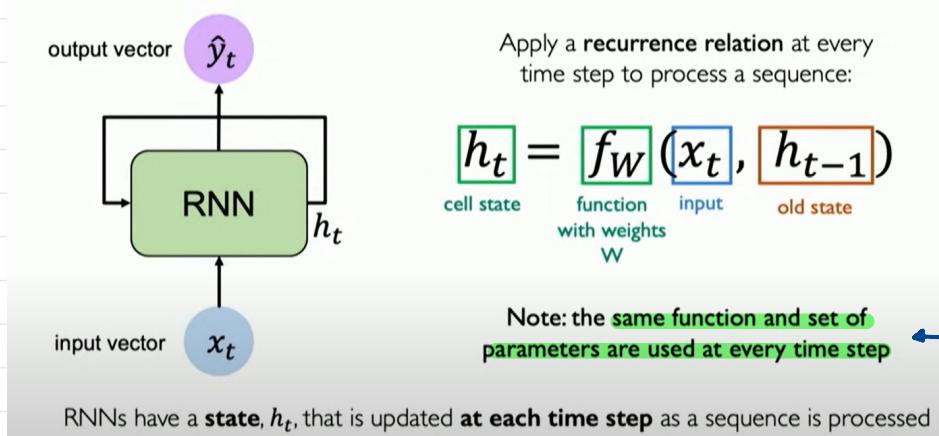


$$\hat{y}_t = f(x_t, h_{t-1})$$

t: time, h_0 : Inner state / memory term

- \hat{y}_t 는 그 때의 입력과 과거 메모리를 모두 담아낸다.

Recurrent Neural Networks (RNNs)



Output Vector

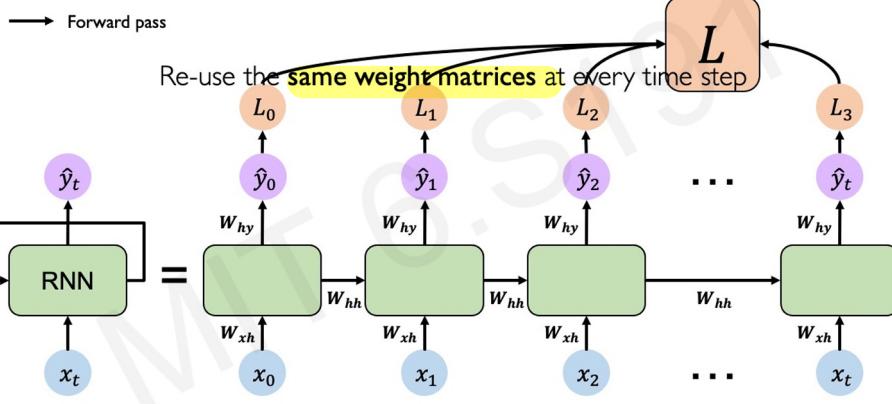
$$\hat{y}_t = W^T h_t$$

Update Hidden State

$$h_t = \tanh(W_h^T h_{t-1} + W_x^T x_t)$$

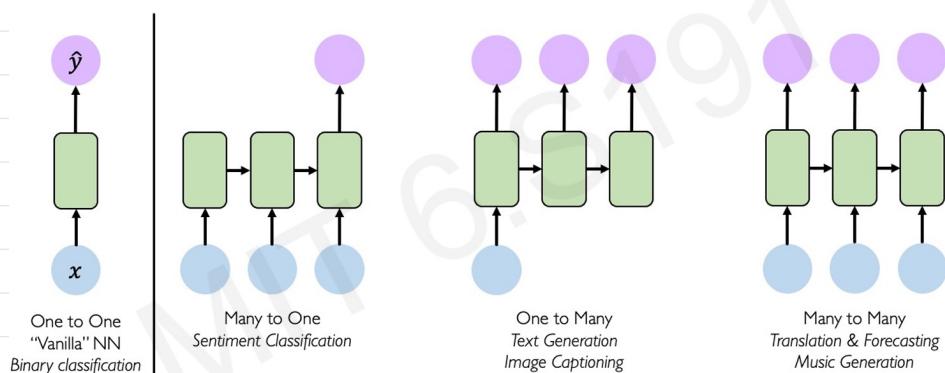
(When input vector is x_t)

W, H_t 은 개별로 초기화된다.



L_t : 각步별 time步의 일시적인 손실
 L : L_t (모든 개별 손실)을 합해
 가중치를 학습하여 총 손실을 구한다.
 $L = \sum_t L_t(\hat{y}_t, y_t)$

RNN Models



What Sequential Models need to deal with

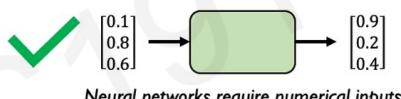
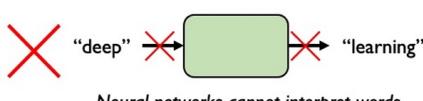
1. Variable-length sequences
2. Track long-term dependencies
3. Maintain information about order
4. Share Parameters across the sequence

→ RNN은 4가지 솔루션 필요하다.

Example of RNN: Word prediction

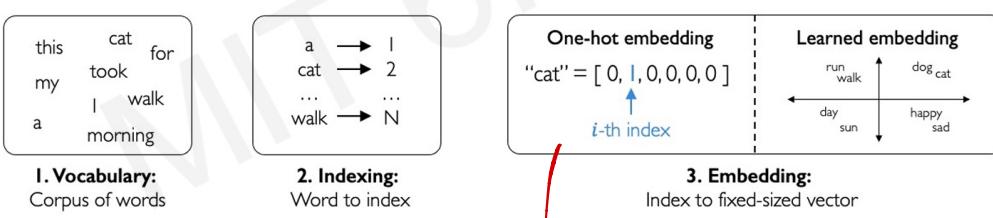
Neural Network은 Text data가 무언가 의미하는지 정확히 알지 못하고 숫자 계산만 할 수 있다:

→ 우리 Text는 Numerical Inputs로 변환해야 한다.



- Training을 통해
 나자신이 의미를 해석하고
 초대한 유사한 숫자 벡터를
 가리키는 것이 목표이다.

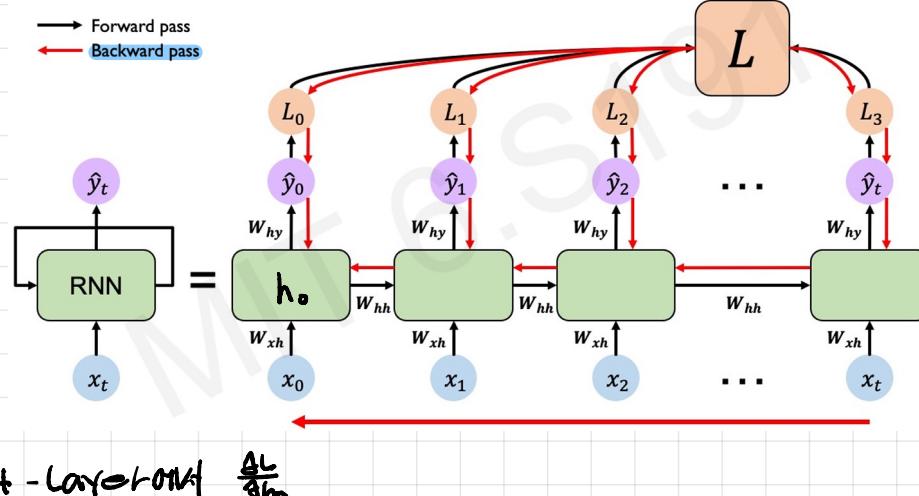
Embedding: transform indexes into a vector of fixed size.



Vector는 단 하나의 값만 1이고 나머지 모든 것은 0인 Encoding

Training Weights of RNN

- Backpropagation Through Time (BPTT)



Ex) 4-Layer RNN

완전한 식

$$\frac{\partial L}{\partial h_0} = \delta_3 W_{hy} f'(h_3) W_{hh} f'(h_2) W_{hh} f'(h_1) W_{hh} + \delta_2 W_{hy} f'(h_2) W_{hh} f'(h_1) W_{hh} + \delta_1 W_{hy} f'(h_1) W_{hh} + \delta_0 W_{hy}$$

이제 줄임 없이 완전한 역전파 과정을 확인할 수 있어! 🎉

→ W_{hh} 가 매우 작다면 $\nabla L \approx 0$ 이 될 수도 있다.

Gradient of h_0 는 유체법 W_{hh} + repeated gradient computation을 포함한다.

Case 1) weight가 높을 때

Exploding gradient 불안정 가능

→ Gradient Clipping으로 해결

Case 2) weight가 매우 낮을 때

(W_{hh})

Vanishing gradient 불안정 가능

Solution 1) Activation function
2) weight initialization
3) Network architecture

→ Gradient가 0인 순간에 지나도
큰 흔적이 아니어지 않도록 한다.

→ long-term data에 대해서 매우 작은
값이 계속 전파되면 오류가 누적된다.

Solution of Vanishing Gradient

1. Activation function

- ReLU를 사용



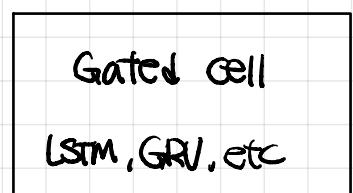
- ReLU는 미분하면 $x > 0$ 에 대해서
항상 1으로 되어야 한다.

2. Parameter Initialize

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

Identity matrix로 weight를
초기화하면 흐름이 된다.

3. Gated Cells



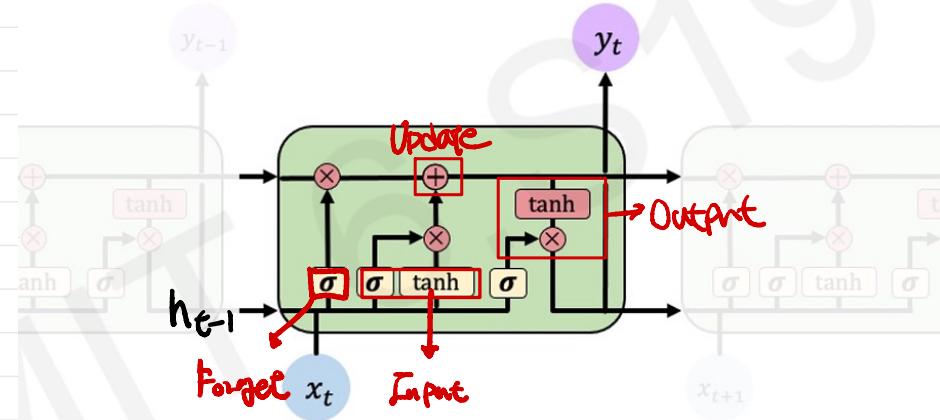
선택적으로 정보를 add or remove
할 수 있도록 한다.
- 중요한 것은 add, 둘 중에서 혹은
remove 한다.

Example of Gated cell : LSTM

LSTM: ANN의 vanishing gradient 문제 해결 가능

Gated LSTM cells control information flow:

1) Forget 2) Store 3) Update 4) Output



→ Uninterrupted Gradient descent을 이용해 Backpropagation을 통해
Vanishing Gradient 문제를 헤쳐나온다.

Limitations of RNNs

1. Encoding bottleneck

- 정보가 압축되면 Encoding을 통해 놓아지고, 이 모든 정보가 흐트러지거나 유실되는 경우가 많다.

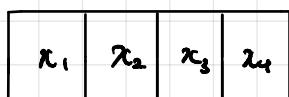
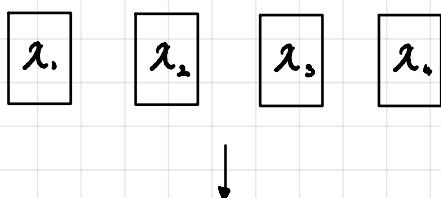
2. 시간별로 처리하면 매우 느릴 수 있다.

→ 병렬 계산을 하는 쉬운 방법도 없다.

3. RNN의 용량이 크지 않아 Long memory를 침식한다.

RNN의 문제점은 극복하려면?

→ 시간 단위의 입력 대신, 한 번에 여러 입력해야 한다.



하지만 RNN을 이용한 경우,
증거는 정보가 사라질 수 있고, data의
증거가 보존되지 않으며, Long term
memory가 침식하다.

Attention : Transformer 모델의 핵심

Ex) 1. 내가 원하는 것 (query)과 제일相近한 것 (key)의 유사도 찾기

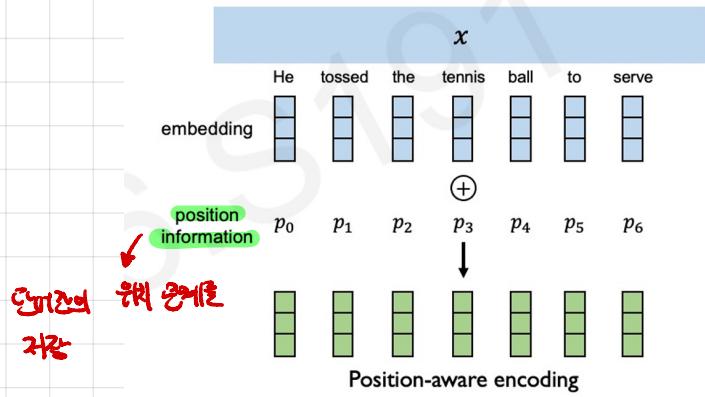
- Attention mask

2. 원하는 것에 기반한 Attention

- Extract values based on attention

Learning Self-Attention With Neural Networks

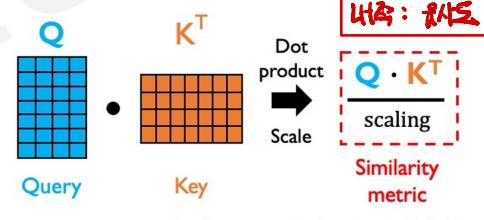
1. Encode position information



3. Compute attention weighting

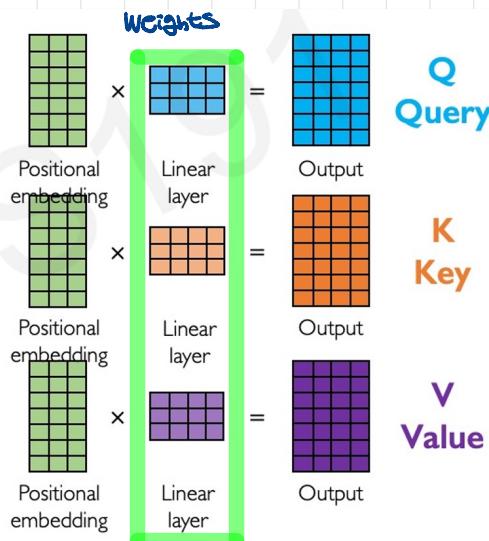
Attention score: compute pairwise similarity between each query and key

How to compute similarity between two sets of features?

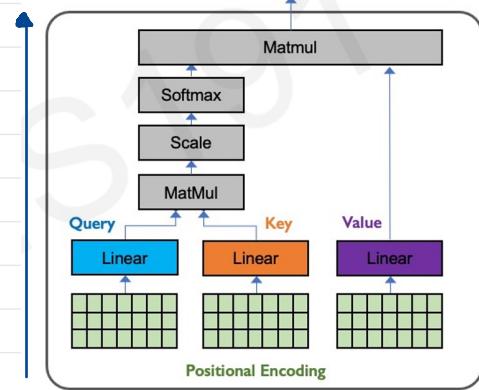


Also known as the "cosine similarity"

2. Extract query, key, value



$$\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right) \cdot V = A(Q, K, V)$$



$$\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right) \cdot V$$

$0 \leq \cdot \leq 1$
→ 높은 중요도를 갖는다.

위의 3가지 과정을 다 거쳐서 마지막으로 Attention을 통해 features를 선택한다.

Self Attention head

각 head는 W_Q, W_K, W_V 을 가지며 각각 다른 블록이나 블록에서 Attention 계산할 수 있다.



Attention weighting



Value



Output



Output of attention head 1



Output of attention head 2



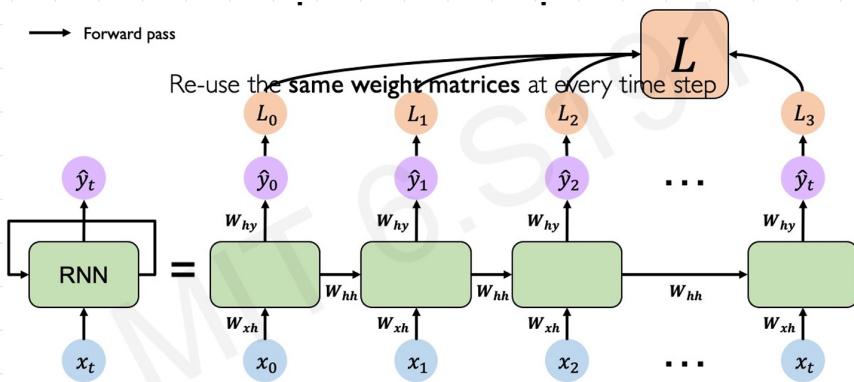
Output of attention head 3

→ 각각의 attention head가 각각의 앤솔, 배경, 유크션을 주목하고 있다.

Multi Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\underbrace{\text{head}_1, \dots, \text{head}_h}_{\text{head를 } h\text{개}}, \underbrace{W_o}_{\text{결과값}})$$

\Rightarrow Long dependency는 허용되는 한정한 수 있다.



$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\in [-1, 1]$$

L_t : Class entropy
 $\text{MSE} \equiv \dots$

- $h_0 = f(x_0 W_{xh})$, $\hat{y}_0 = W_{hy} \cdot h_0$, $L_0 = Y_0 - \hat{y}_0$
- $h_1 = f(h_0 W_{hh} + x_1 W_{xh})$, $\hat{y}_1 = h_1 \cdot W_{hy}$, $L_1 = Y_1 - \hat{y}_1$
- $h_2 = f(h_1 W_{hh} + x_2 W_{xh})$, $\hat{y}_2 = h_2 \cdot W_{hy}$, $L_2 = Y_2 - \hat{y}_2$
- $h_3 = f(h_2 W_{hh} + x_3 W_{xh})$, $\hat{y}_3 = h_3 \cdot W_{hy}$, $L_3 = Y_3 - \hat{y}_3$
- $L = \sum_{t=0}^3 L_t$

Gradient - Multi Layer ANN

$$\begin{aligned} \frac{\partial L}{\partial W_{hy}^{(4)}} &= \sum_{t=0}^T \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial W_{hy}^{(4)}} \\ \frac{\partial L}{\partial W_{hh}^{(1)}} &= \sum_{t=0}^T \frac{\partial L_t}{\partial h_t^{(1)}} \cdot \frac{\partial h_t^{(1)}}{\partial W_{hh}^{(1)}} \\ \frac{\partial L}{\partial W_{xh}^{(1)}} &= \sum_{t=0}^T \frac{\partial L_t}{\partial h_t^{(1)}} \cdot \frac{\partial h_t^{(1)}}{\partial W_{xh}^{(1)}} \end{aligned}$$

Example of RNN (4 layers) - Feed forward

Data 2: $\begin{bmatrix} 0.5 & 0.1 & -0.3 \\ -0.2 & 0.4 & 0.3 \\ 0.1 & -0.5 & 0.2 \\ 0.4 & 0.2 & -0.1 \end{bmatrix}$ $Y: \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$ $W_{xh} = \begin{bmatrix} 0.2 & -0.3 \\ 0.4 & 0.1 \\ -0.5 & 0.2 \end{bmatrix}$ $W_{hh} = \begin{bmatrix} 0.6 & -0.2 \\ 0.1 & 0.5 \end{bmatrix}$ $W_{hy} = \begin{bmatrix} 0.3 & -0.7 \\ -0.5 & 0.8 \end{bmatrix}$

$$h_0 = f(W_{xh} \cdot x_0) = \begin{bmatrix} 0.2 & -0.3 \\ 0.4 & 0.1 \\ -0.5 & 0.2 \end{bmatrix}^T \cdot \begin{bmatrix} 0.5 \\ 0.1 \\ 0.3 \end{bmatrix} = \begin{bmatrix} (0.2 \times 0.5) + (0.4 \times 0.1) + (-0.5 \times 0.3) \\ (-0.3 \times 0.5) + (0.1 \times 0.1) + (0.2 \times 0.3) \end{bmatrix} = \tanh \begin{bmatrix} 0.28 \\ -0.198 \end{bmatrix} = \begin{bmatrix} 0.278 \\ -0.198 \end{bmatrix}$$

$$\hat{y}_0 = W_{hy} \cdot h_0 = \begin{bmatrix} 0.3 & -0.7 \\ -0.5 & 0.8 \end{bmatrix} \cdot \begin{bmatrix} 0.278 \\ -0.198 \end{bmatrix} = \begin{bmatrix} 0.2205 \\ -0.2947 \end{bmatrix}$$

* h_t, \hat{y}_t, L_t 는 모두 다음과 같다:
 (W_{hh}, W_{xh}, W_{hy}) 를 초기화.

$$L_0 = \frac{1}{2} \|Y_0 - \hat{y}_0\|^2 = \frac{1}{2} [(1 - 0.2205)^2 + (0 - (-0.2947))^2] = 0.3478$$

$$\text{Final } L: \sum_{t=0}^T L_t$$

BPTT formula induction

$$\frac{\partial L}{\partial W_{hy}} = \frac{1}{2} \frac{\partial}{\partial W_{hy}} \left(\sum_{t=0}^T L_t \right) = \frac{1}{2} \frac{\partial}{\partial W_{hy}} \left(\sum_{t=0}^T \|W_{hy} h_t - Y_t\|^2 \right) = \sum_{t=0}^T (\|\hat{y}_t - Y_t\|) \cdot h_t^T = \sum_{t=0}^T \delta_t \cdot h_t^T$$

$$\begin{aligned} \frac{\partial L}{\partial W_{hh}} &= \frac{1}{2} \frac{\partial}{\partial W_{hh}} \left(\sum_{t=0}^T L_t \right) = \frac{1}{2} \frac{\partial}{\partial W_{hh}} \left(\sum_{t=0}^T \|W_{hy} h_t - Y_t\|^2 \right) = \frac{1}{2} \frac{\partial}{\partial W_{hh}} \left(\sum_{t=0}^T (\|W_{hy} \cdot \tanh(W_{xh} \cdot x_t + W_{hh} \cdot h_{t-1}) - Y_t\|^2) \right) \\ &= \sum_{t=0}^T \delta_t W_{hy}^T (1 - h_t^2) h_{t-1}^T \quad (1 - \tanh^2(z)) \end{aligned}$$

$$\frac{\partial L}{\partial W_{xh}} = \frac{1}{2} \frac{\partial}{\partial W_{xh}} \left(\sum_{t=0}^T L_t \right) = \frac{1}{2} \frac{\partial}{\partial W_{xh}} \left(\sum_{t=0}^T \|W_{hy} h_t - Y_t\|^2 \right) = \frac{1}{2} \frac{\partial}{\partial W_{xh}} \left(\sum_{t=0}^T (\|W_{hy} \cdot \tanh(W_{xh} \cdot x_t + W_{hh} \cdot h_{t-1}) - Y_t\|^2) \right)$$

$$\cdot \sum_{t=0}^T \delta_t W_{hy}^T (1 - h_t^2) x_t^T$$