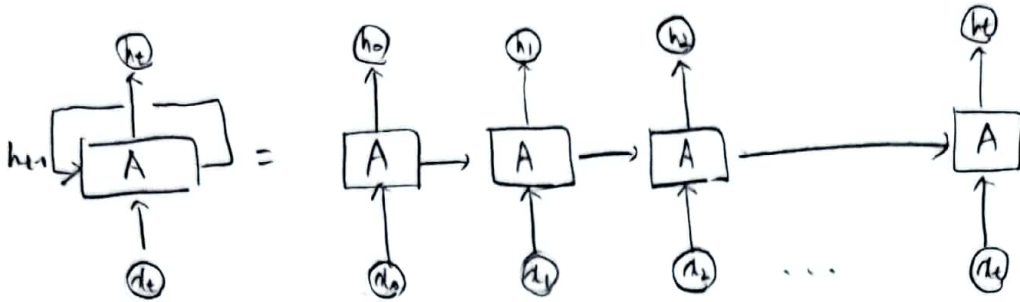


# RNN (Recurrent Neural Network)

Sequence가 입력되는 순서대로 주어질 때,

각 time step에서 들어오는 입력 벡터  $d_t$ 과 그 전 time step의 RNN으로부터 계산한 hidden state 벡터  $h_{t-1}$ 을 입력으로 받아 현재 time step에서의  $h_t$ 을 출력으로 내어주는 구조



여기 time step별 동일한 파라미터를 재사용 가능  
(A)

hidden state 계산 방법

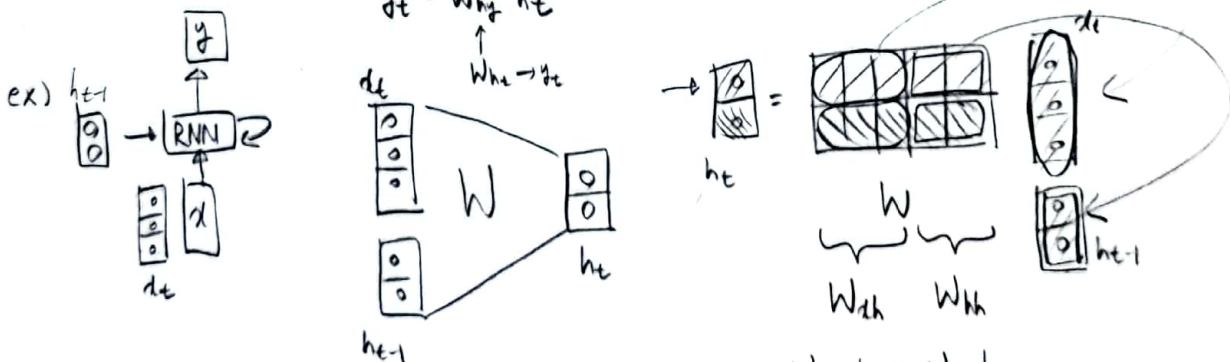
- $h_{t-1}$ : old hidden-state vector
- $d_t$ : input vector at same time step
- $h_t$ : new hidden-state vector
- $f_w$ : RNN function with parameters  $w$
- $y_t$ : output vector at time step  $t \rightarrow$  이 time step에서 계산하는 경우와  
이전 time step에서 계산하는 경우로 나뉨

$$h_t = f_w(h_{t-1}, d_t)$$

$\rightarrow$  이 time step에서 계산

$$h_t = f_w(h_{t-1}, d_t) \rightarrow h_t = \tanh(W_{hh}h_{t-1} + W_{hd}d_t)$$

$y_t = W_{hy}h_t$



$y_t$ : 이진출력: scalar  $\rightarrow$  sigmoid  
다중출력: vector  $\rightarrow$  softmax  
dim: #class

$$= W_{hd}d_t + W_{hh}h_{t-1}$$

$W_{hd} \rightarrow d_t \rightarrow h_t$        $W_{hh} \rightarrow h_{t-1} \rightarrow h_t$

# RNN 유형

## ① One-to-One

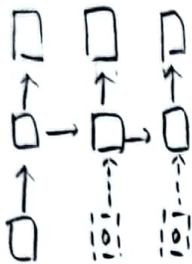


ex) [키, 몸무게, 나이] → 리얼업, 2판업!

standard neural networks

입력, 출력 모두 sequence가 아님 (1개의 time step)

## ② One-to-many



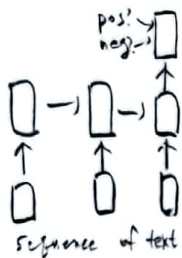
입력이 1개의 time step이고 출력은 여러 time step

ex) 이미지 캡셔닝

하나의 이미지를 입력으로 주면 이미지에 대한 설명을 생성

입력: 첫번째 time step에만 들어감 (나머지 입력은 0으로 간주)

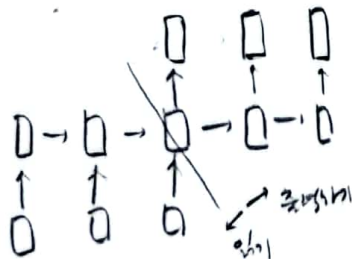
## ③ Many-to-one



입력은 sequence로 받아 여러 time step에서 출력

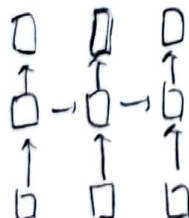
ex) 감성 분석

## ④ Many-to-many (Sequence-to-sequence)



입력, 출력 모두 sequence

ex) 기계 번역



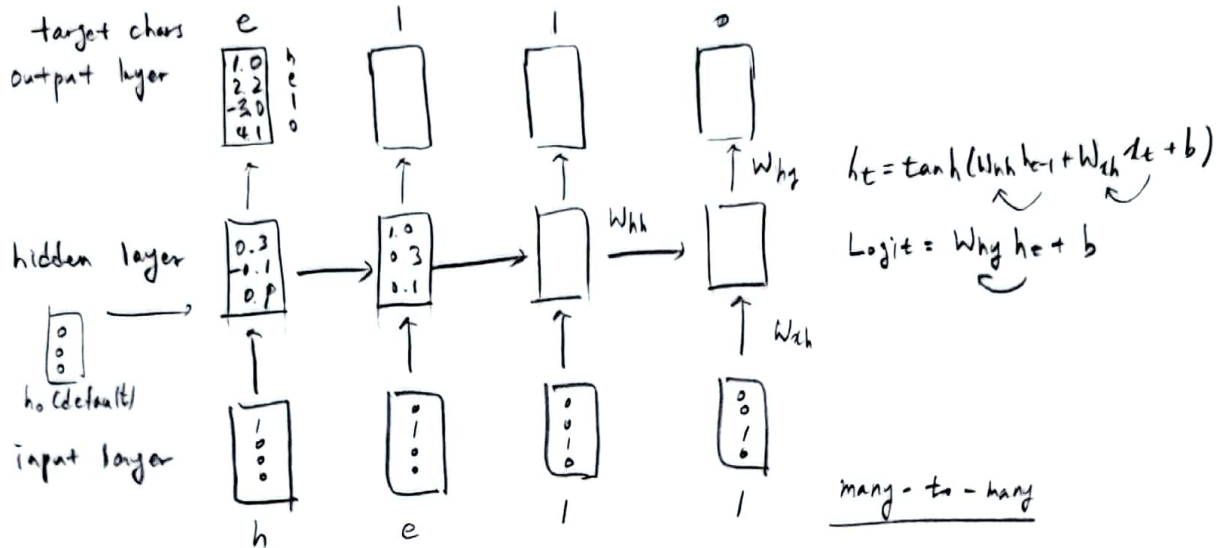
ex) Video Classification on frame level  
pos tagging.

# Character-level language model

언어 모델 이전에 등장한 문자열 생성 기법으로 다음 단어를 예측하는 task  
여기서는 문자 단위.

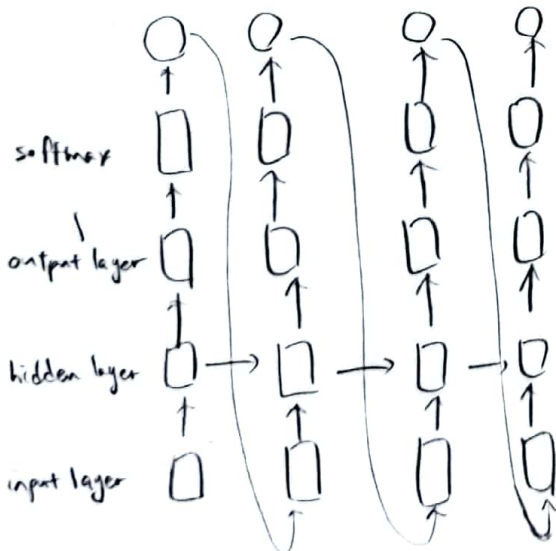
ex) "hello" 가 주어졌을 때 다음 단어를

Vocabulary: {'h', 'e', 'l', 'o'},  $h = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$   $e = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$



h가 주어지면 e, e가 주어지면 l, ... 을 예측

test time 예시

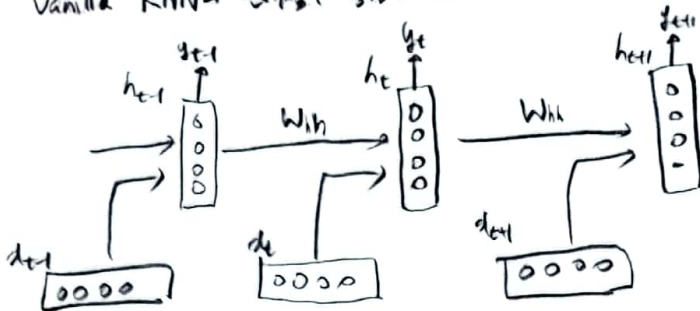


ex) 전 time step 가니 주어진 것을 바탕으로  
다음은 주어진 예측,  
세익스피어의 죽음 (주안),  
이웃방 여사,  
latex 논문,  
C 프로그래밍.

# BPTT (Backpropagation through time)

Loss function에 대한 backpropagation은 동시  $W_{xh}$ ,  $W_{hh}$ ,  $W_{hy}$ 를 학습  
 sequence가 길어지면 전체 sequence의 loss를 계산하고 gradient를 계산하는 것은 한계가 존재  
 Truncation: 순이론에 따라 계산의 길이로 학습  
 특징: hidden state를 시작하자면 time step이 진행됨에 따라 중요한 부분은 잘 학습하는 안 수 있다.  
 ex) Zf statement cell, route detection cell...

Vanilla RNN의 문제점: 동일한 matrix를 여러 time step마다 곱하면서



동일한 matrix를 backpropagation에서도 마찬가지로 곱하는 수가 계속 증가하므로

Vanishing Gradient 또는 Exploding Gradient가 발생  
 (값이 < 1) (값이 > 1)

예시)  $h_t = \tanh(W_{xh}x_t + W_{hh}h_{t-1} + b)$ ,  $t=1, 2, 3$

For  $W_{xh}=3$ ,  $W_{hh}=2$ ,  $b=1$

$$h_3 = \tanh(2x_3 + 3h_2 + 1)$$

$$h_2 = \tanh(2x_2 + 3h_1 + 1)$$

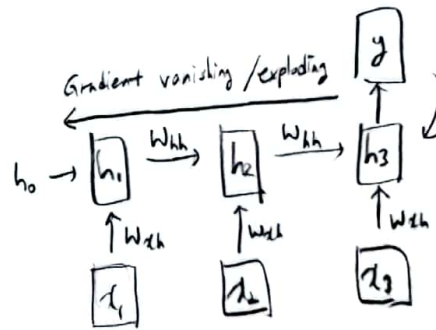
$$h_1 = \tanh(2x_1 + 3h_0 + 1)$$

$$h_3 = \tanh(2x_3 + 3 \tanh(2x_2 + 3 \tanh(2x_1 + 3h_0 + 1) + 1) + 1)$$

301 time step마다 계속 곱함 → exploding

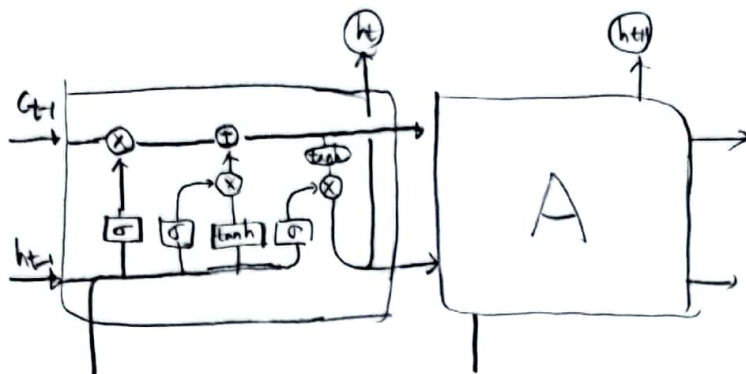
만약  $W_{hh}=0.2$ 이면 → vanishing

(Long-term dependency 문제)



## LSTM (Long Short-Term Memory)

RNN에서 hidden state vector를 단기 기억으로 보면 time step이 진행됨에 따라 단기 기억은 보다 길게 기억할 수 있도록 개선할 노면



$$h_t = f_w(x_t, h_{t-1})$$

(RNN)

$$\{c_t, h_t\} = \text{LSTM}(x_t, c_{t-1}, h_{t-1})$$

↑                  ↑                                  ↑                  ↑

중어  
순환형  
정보 흐름

세포 상태를 저장  
해서 다음 단계에서  
사용하는 역할이 있음.  
리본만 나열함

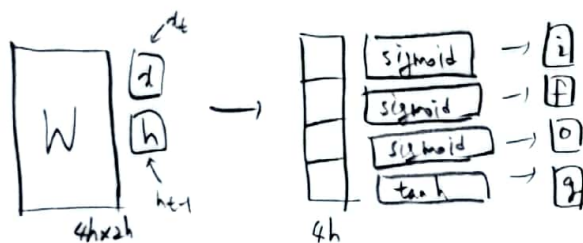
cell state      hidden state

- $i$ : input gate, cell에 쓸지 말지 결정,  $\sigma$ 를 통해 0~1 범위의 값, cell/hidden state 두 입력으로 이루어짐
- $f$ : forget gate, 정보를 어느 정도 지우는지,  $\sigma$ 를 통해 0~1 범위의 값
- $O$ : output gate, cell 정보를 어느 정도 hidden state에서 사용하냐,  $\sigma$ 를 통해 0~1 범위의 값
- $g$ : Gate gate, 어느 정도로 cell state에 반영하냐,  $\tanh$ 를 통해 -1~1 범위의 값

$f$ : forget gate. 정보를 어느 정도 지우는지, 0을 통해 0~1 범위의 값

O: output gate, cell 정보를 어느 정도 hidden state에서 내놓을지, 0을 통해 0~1 범위 내 값

g. Gate gate, 어느정도 cell state를 받아들이지, tanh를 통해 -1~1 범위의 값



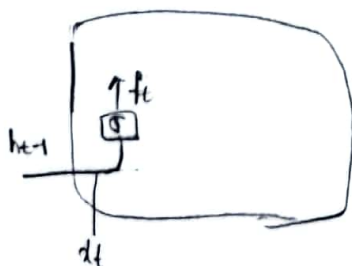
$$\begin{pmatrix} \dot{z} \\ \dot{f} \\ \dot{o} \\ \dot{g} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ \text{dt} \end{pmatrix}$$

$$C_t = f \odot C_{t-1} + i \odot g$$

$$h_t = 0 \odot \tanh(C_t)$$

IFOG는 전 time에서 얻은  $C_{n-1}$ 은 각각씩에 변함

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



ex)  $G_{t-1} = \begin{bmatrix} 3 \\ 5 \\ -2 \end{bmatrix}$

$$f_t = \begin{bmatrix} 0.7 \\ 0.4 \\ 0.8 \end{bmatrix} \quad 0.1 = 1$$

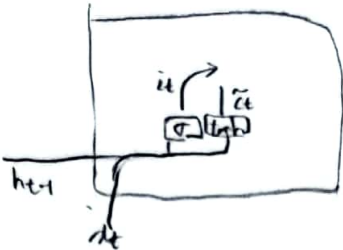
element-wise product:  $\text{sum} \begin{bmatrix} 2.1 \\ 2.0 \\ -1.6 \end{bmatrix}$  or  $\text{sum}$



Input gate:  $\hat{z}_t = \sigma(W_i \cdot [h_{t-1}, d_t] + b_i)$

Gate gate:  $\tilde{C}_t = \tanh(W_c \cdot [h_{t-1}, d_t] + b_c)$

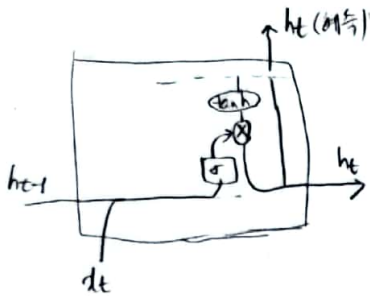
new cell state:  $C_t = \hat{f}_t \cdot C_{t-1} + \hat{z}_t \cdot \tilde{C}_t$



한 번의 선택된 값으로  $C_t$ 에 어떠한 정보를 넣을지 여부를  
그보다 좀 더 큰 값인  $\tilde{C}_t$ 로 만들려는 것.  $\hat{z}_t$ 로 이용해 인박만  
발생시켜, 실제로 제거하고 싶은 정보는 2단계에 걸쳐서 나온다

Output gate:  $O_t = \sigma(W_o \cdot [h_{t-1}, d_t] + b_o)$

hidden state:  $h_t = O_t \cdot \tanh(C_t)$



$C_t$ 와  $h_{t-1}$  차이?

$\rightarrow C_t$ : 기억해야 하는 모든 정보를 담은 벡터.

$h_t$ : 현재 time step에서 output layer에 입력으로 사용되는 벡터,  
다음 time step에 예측하기 위해 필요한 정보를 담음(필터링)

ex) "hell" 아냐

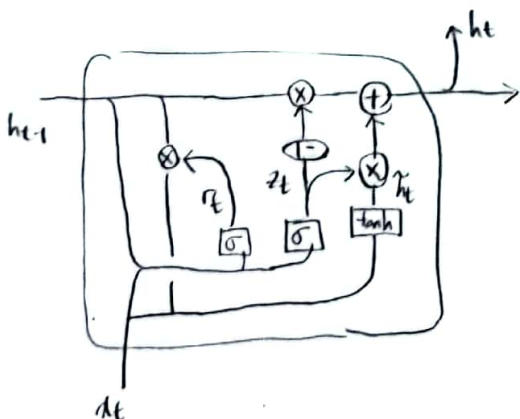
다음어가 띄어쓰는 것은 항상 필요한 정보는 아니지만  
기억해 두는 것이 좋다. 해당 정보를  $C_t$ 에 저장  
다음에 아가 되나 할 때는 항상 다음어가 필요하기 때문  
필터링한 부분은  $h_t$ 에 저장

## GRU (Gated Recurrent Unit)

LSTM의 경량화, 2개의 게이트, 1개의 hidden state

cell state와 hidden state를 hidden state로 전환

이때의 hidden state는 LSTM의 cell state와 비슷한 것



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

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

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

$$h_t = (1 - z_t) \cdot h_{t-1} + z_t \cdot \tilde{h}_t \quad (\text{가중치})$$

$$C_t = \hat{f}_t \cdot C_{t-1} + \hat{z}_t \cdot \tilde{C}_t$$

$$\text{ex) } z_t = \begin{bmatrix} 0.6 \\ 0.3 \\ 0.5 \end{bmatrix}$$

$$1 - z_t = \begin{bmatrix} 0.4 \\ 0.7 \\ 0.5 \end{bmatrix}$$

$\rightarrow$  가중치는 gradient back: long-term dependency 비관.