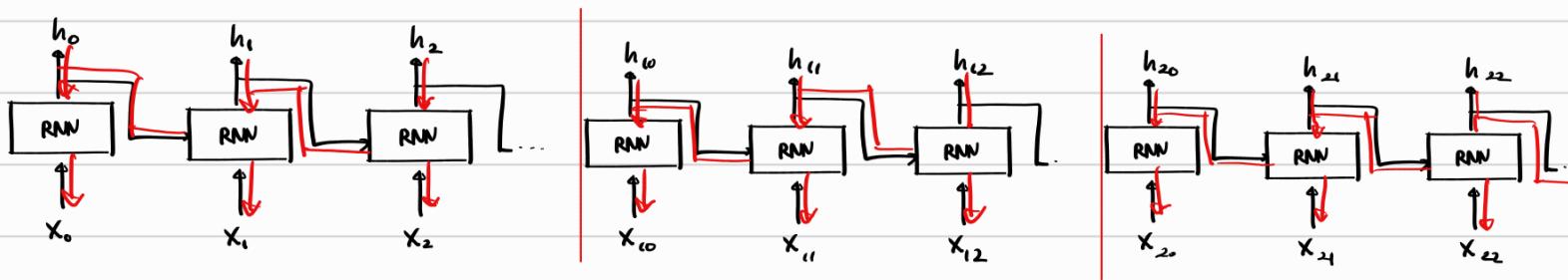
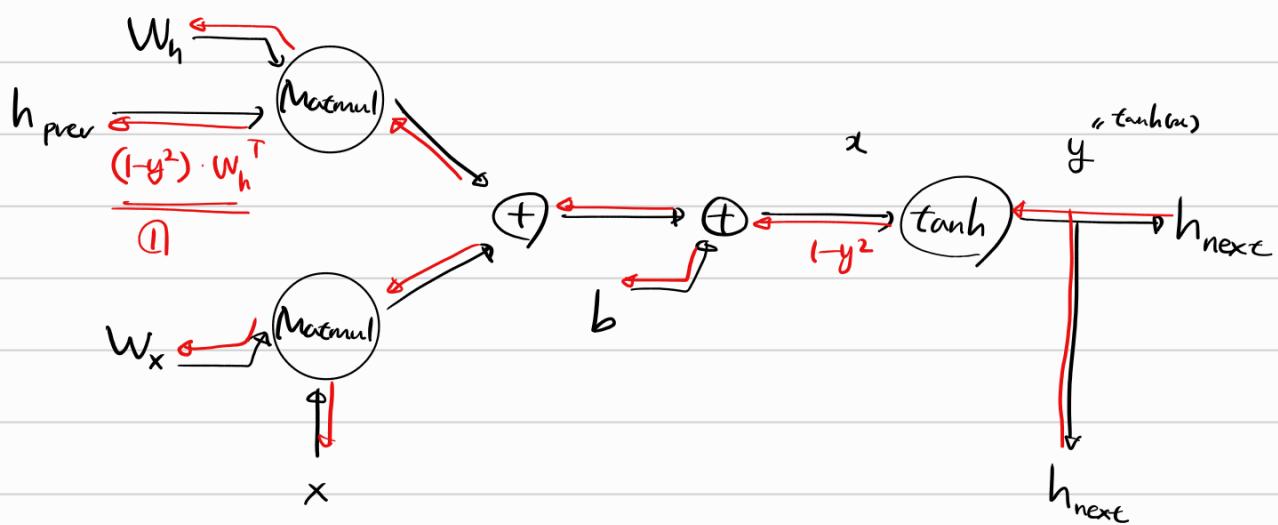
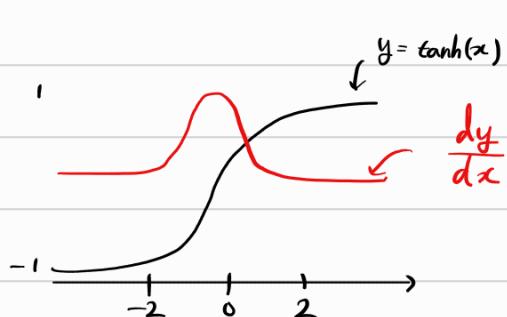


CH. 1 RNN의 한계

Recall. RNN



Thm 1.1 gradient exploding & Vanish on RNN



P1. $dh_{prev} = (1-y^2) \times W_h^T$ 만약 W_h 의 특이값이 보다 크면?
or $y=0$?
could admit Nan

P2. $1-y^2$ 은 monotony decrease.

그리고 W_h^T 도 역시 소실을 야기할 수 있다.

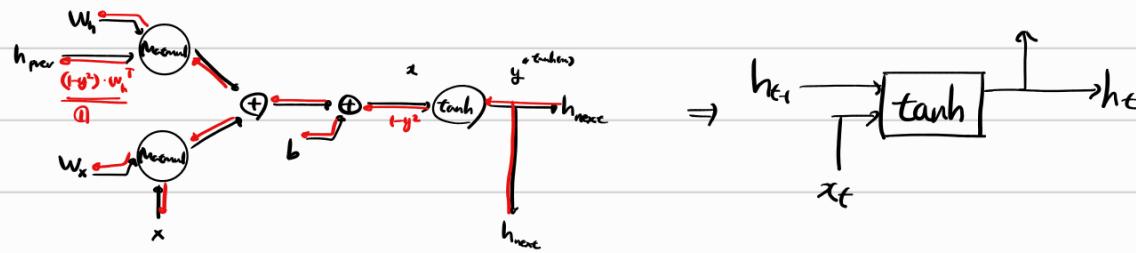
That is, dh_{prev} 가 소실, 폭발 가능성이 있다. 이는 장기 기억력에 악영향을 준다!

Thm 1.2 gradients clipping can help gradient exploding

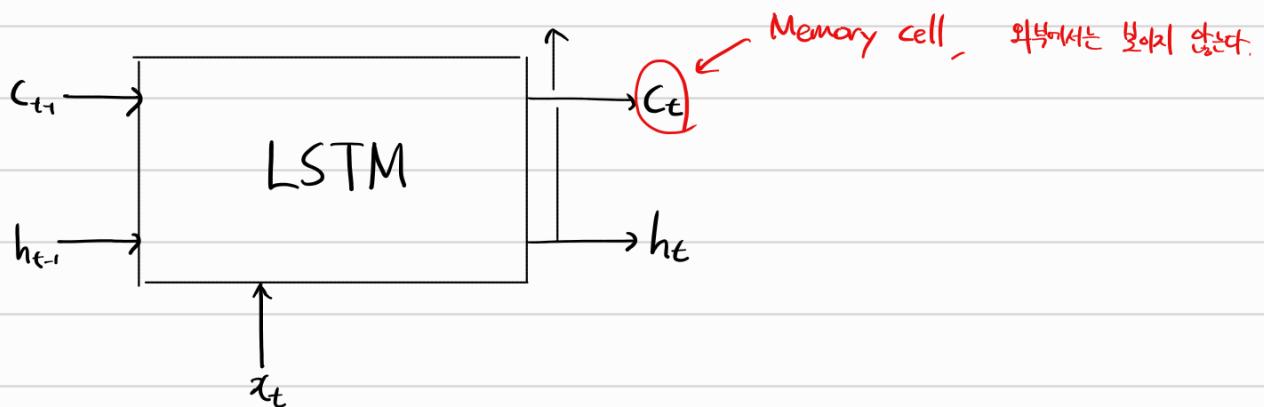
If $\|\hat{g}\| \geq \text{threshold}$, Then $\hat{g}' = \frac{\text{Threshold}}{\|\hat{g}\|} \times \hat{g}$

CH.2 gradients vanish and LSTM

Notation



Def 2.1 LSTM Architecture



Def 2.2 output gate

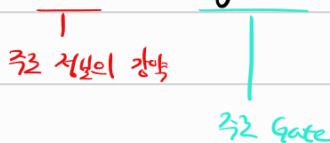
Let $O = \sigma(x_t w_x^{(0)} + h_{t-1} w_h^{(0)} + b^{(0)})$

↑ Sigmoid ↑ 0~1, 값이 클수록 확보 비중 ↑

Then $h_t = O \odot \tanh(c_t)$

↑ hadamard Product

Note 2.3 tanh vs sigmoid



Def 2.3 forget gate

Let $f = \sigma(x_t w_x^{(f)} + h_{t-1} w_h^{(f)} + b^{(f)})$

Then $c_t = f \odot \tanh(c_{t-1})$

Def 2.4 New memory cell

Let $\tilde{g} = \tanh(x_t w_x^{(g)} + h_{t-1} w_h^{(g)} + b^{(g)})$

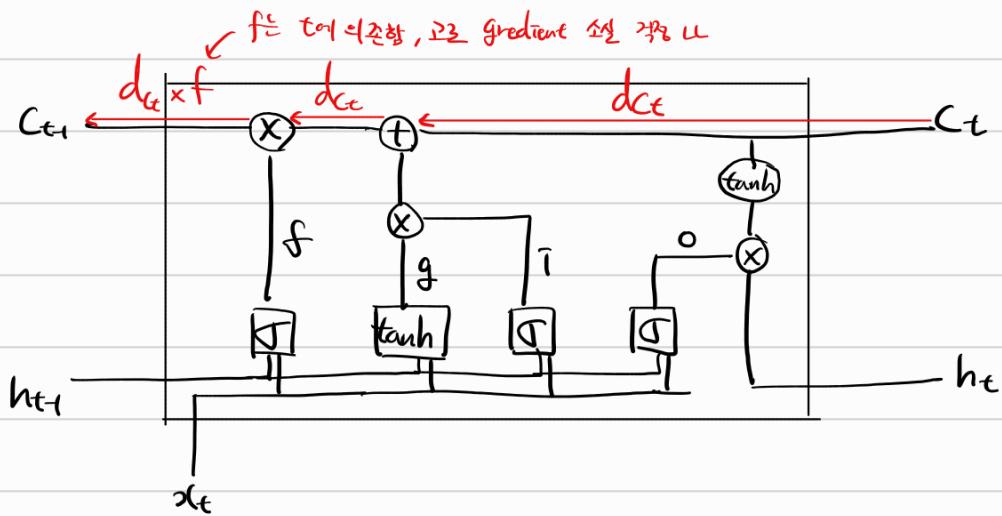
\tilde{g} 기억 셀에 더해진 새로운 기억

Def 2.5 Input gate

Let $\tilde{i} = \sigma(x_t w_x^{(i)} + h_{t-1} w_h^{(i)} + b^{(i)})$

Then $C_t += \tilde{i} \odot g$

Rmk 2.6 Memory cell's BP on LSTM



CH.3 LSTM의 구현과 개선

$$f = \sigma(h_{t-1} W_h^{(f)} + x_t W_x^{(f)} + b^{(f)})$$

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

$$\bar{t} = \sigma(h_{t-1} W_h^{(\bar{t})} + x_t W_x^{(\bar{t})} + b^{(\bar{t})})$$

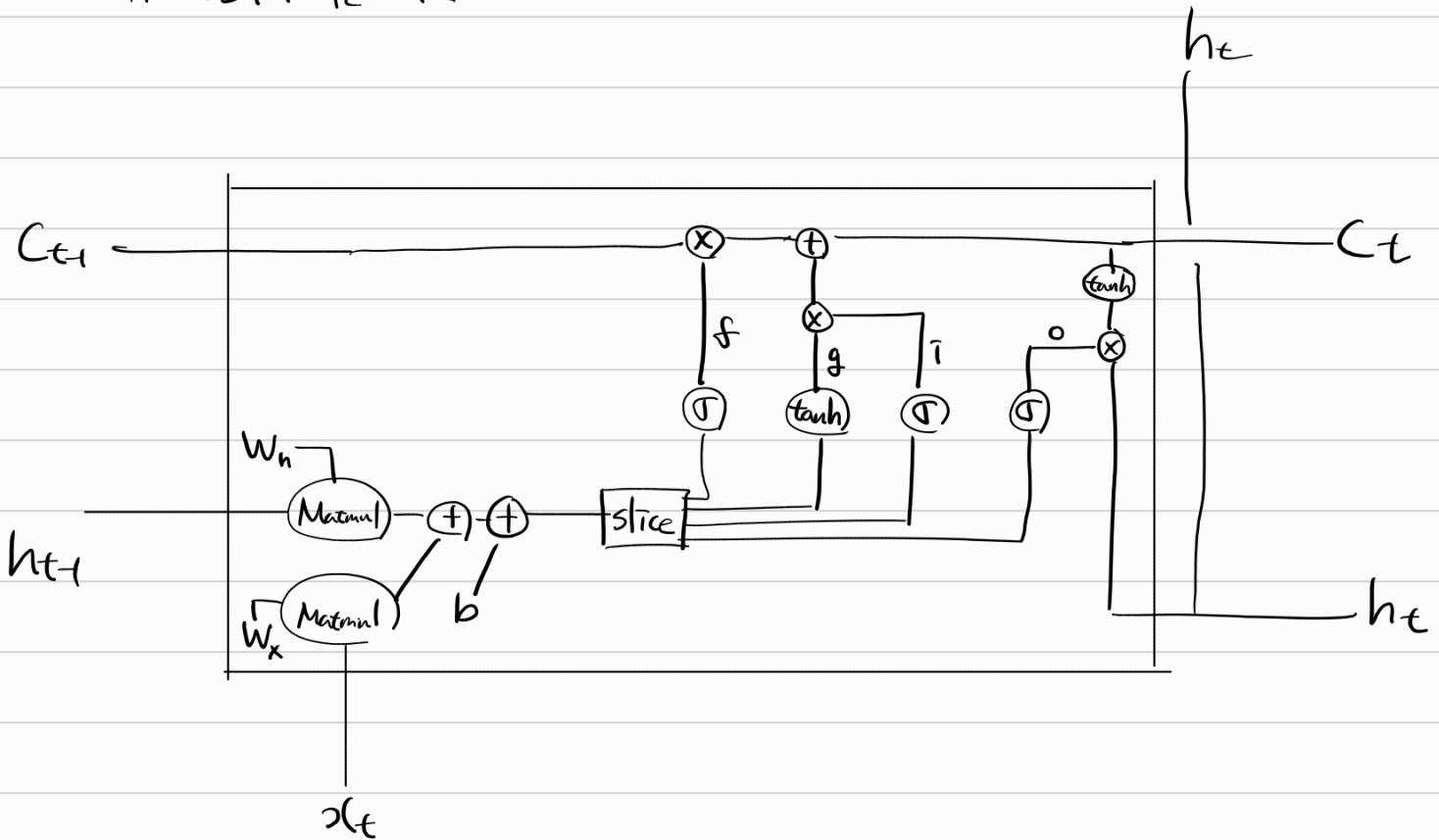
$$o = \sigma(h_{t-1} W_h^{(o)} + x_t W_x^{(o)} + b^{(o)})$$

$$C_t = f \circ C_{t-1} + \bar{t} \circ g$$

$$h_t = \tanh(C_t)$$

$$\Rightarrow x_t \underbrace{[W_x^{(f)} W_x^{(g)} W_x^{(\bar{t})} W_x^{(o)}]}_{W_x} + h_{t-1} \underbrace{[W_h^{(f)} W_h^{(g)} W_h^{(\bar{t})} W_h^{(o)}]}_{W_h} + \underbrace{[b^{(f)} b^{(g)} b^{(\bar{t})} b^{(o)}]}_b$$

Note 3.1 LSTM 계산 과정

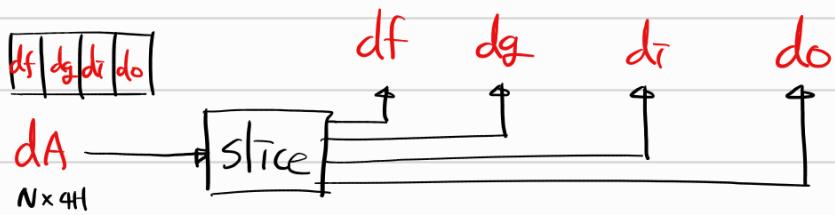
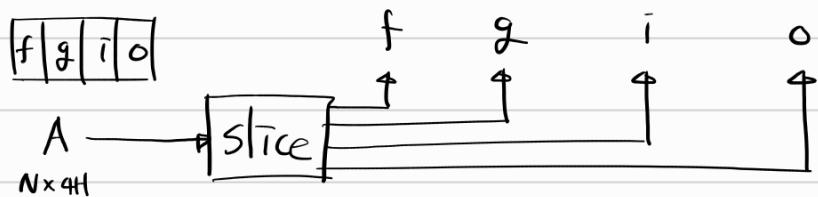


Thm 3.2 matrix size of LSTM

$$X_t W_X + h_{t-1} W_h + b$$

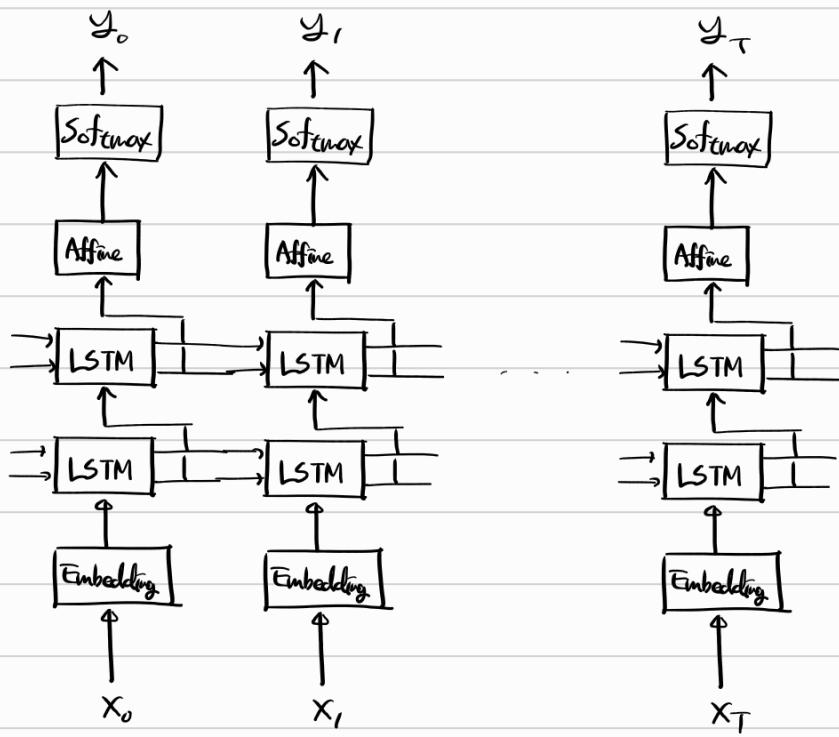
$N \times D$ $D \times 4H$ $N \times H$ $H \times 4H$ $N \times H$

Thm 3.3 BP of LSTM slice



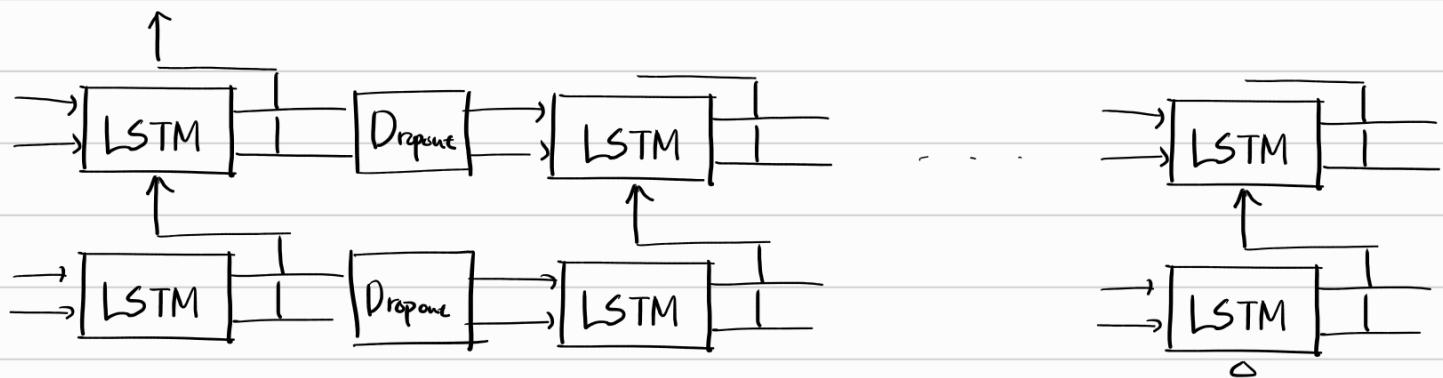
$\Rightarrow \text{np.hstack}(df, dg, di, do)$

Def. 3.4 LSTM based RNNLM

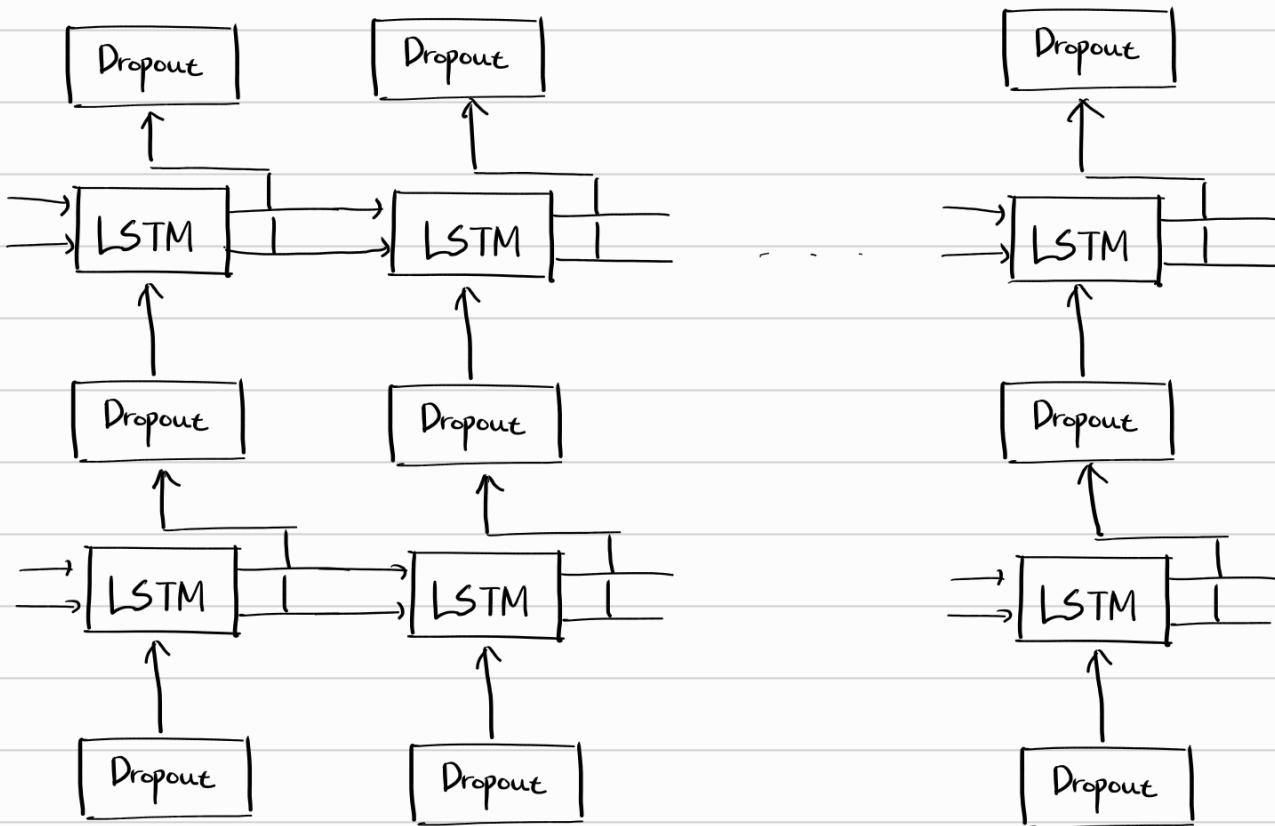


Thm 3.5 Dropout on LSTM

Option ① : 시계열 방향 합입, 이는 학습에서 시간이 지남에 따라 정보가 사라질 수 있다. 즉, 끝은 시간에 비례해서 노이즈 축적



option ② : 상하 방향 합입



=) 시간에 따른 정보 손실 X

option ③ Variational Dropout

option ① + ②에서 같은 포지션 계층끼리 마스크를 공유!

+ Drop Connect라는 기술도 있음

Thm 3.6 weight tying

먼저 매개변수를 줄이면 좋은 것 ① 학습 시간 ↓

② overfitting 방지

