
LSTM and GRU

(Long Short Term Memory and Gate Recurrent Unit)

Winter Vacation Capstone Study

TEAM Kai.Lib

발표자 : 김수환

2020.01.06 (MON)

LSTM의 등장배경

▪ 게이트가 추가된 RNN

RNN은 순환 경로를 포함하여 과거의 정보를 기억할 수 있었다. 구조가 단순하여 구현도 쉽게 할 수 있었지만 안타깝게도 성능이 좋지 못하다. 그 원인은 많은 경우 시계열 데이터에서 시간적으로 많이 떨어진 장기long term 의존 관계를 잘 학습할 수 없다는 데 있다. 해서 요즘에는 단순한 RNN 대신 LSTM이나 GRU라는 계층이 주로 쓰인다.

LSTM이나 GRU에는 게이트gate라는 구조가 더해져 있는데, 이 게이트 덕분에 시계열 데이터의 장기 의존 관계를 학습할 수 있다. 이번 스터디에서는 LSTM과 GRU 내부적으로 어떠한 구조로 되어있으며, 어떻게 이 구조가 '장기 기억'을 가능하게 하는지를 이해한다

Tom was watching TV in his room. Mary came into the room. Mary said hi to ?

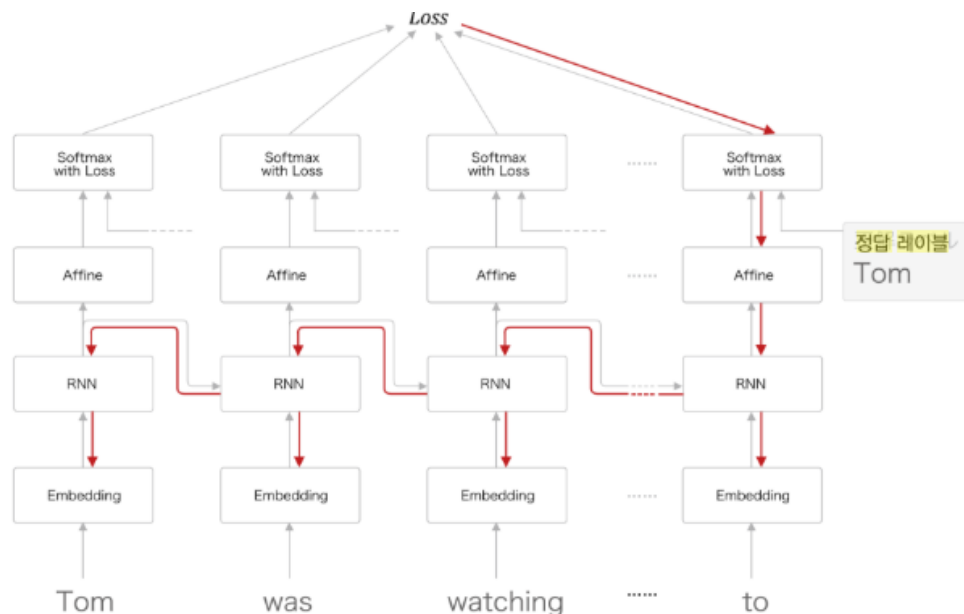
"?"에 들어가는 단어는 "Tom"이다. RNN에서 이 문제에 올바르게 답하려면, 현재 맥락에서 "Tom이 방에서 TV를 보고 있음"과 "그 방에 Mary가 들어옴"이란 정보를 기억해둬야 한다. 다시 말해 이런 정보를 RNN 계층의 은닉 상태에 인코딩해 보관해둬야한다.

LSTM의 등장배경

▪ RNN의 문제점

Tom was watching TV in his room. Mary came into the room. Mary said hi to ?

그럼 이 예를 RNNLM 학습의 관점에서 생각해보면 여기서 정답 레이블로 "Tom"이라는 단어가 주어졌을 때, RNNLM에서 기울기가 어떻게 전파되는지를 살펴보자. 물론 학습은 BPTT로 수행한다. 따라서 정답 레이블이 "Tom"이라고 주어진 시점으로부터 과거 방향으로 기울기를 전달하게 된다. 그림으로는 다음처럼 된다.

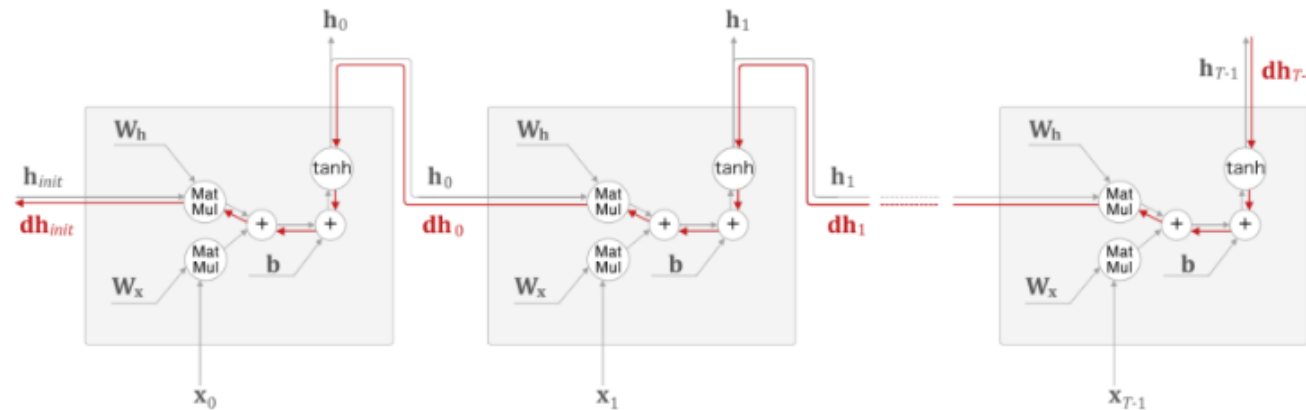


다음 그림처럼 RNN 계층이 과거 방향으로 의미 있는 기울기를 전달함으로써 시간 방향의 의존 관계를 학습할 수 있는 것이다. 하지만 현재의 단순한 RNN 계층에서는 시간을 거슬러 올라갈수록 기울기가 작아지거나 (**기울기 소실**) 혹은 커질 수 있으며 (**기울기 폭발**), 대부분 둘 중 하나의 운명을 걷게 된다.

RNN의 문제점

기울기 소실과 기울기 폭발의 원인

그럼 RNN 계층에서 기울기 소실(혹은 기울기 폭발)이 일어나는 원인을 살펴보자. 생각을 간단하게 하기 위해 다음 그림과 같이 RNN 계층에서의 시간 방향 기울기 전파에만 주목해보자.

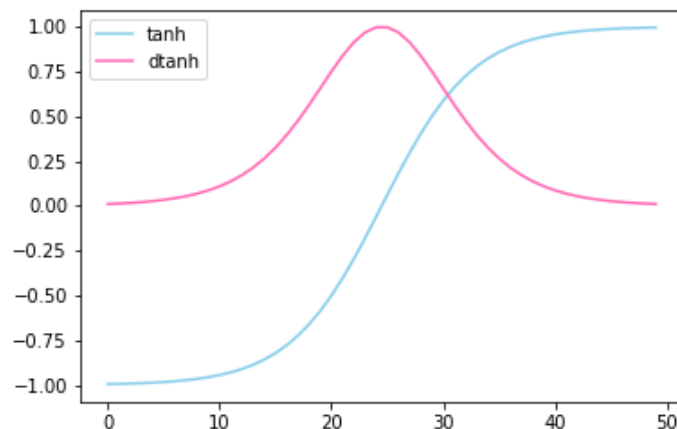


길이가 T 인 시계열 데이터를 가정하여 T 번째 정답 레이블로부터 전해지는 기울기가 어떻게 변하는지 살펴보자. 시간 방향 기울기에 주목하면 역전파로 전해지는 기울기는 차례로 'tanh', '+', 'Matmul(행렬곱)' 연산을 통과한다는 것을 알 수 있다. '+'의 역전파는 상류에서 전해지는 기울기를 그대로 하류로 흘려보낼 뿐이니, 'tanh'와 'Matmul'에만 주목해보자. 우선 'tanh'부터 보자.

RNN의 문제점

▪ 기울기 소실과 기울기 폭발의 원인 - tanh

앞의 스터디에서 tanh 함수의 미분에서 봤듯이, $y = \tanh(x)$ 일 때의 미분은 $\frac{\partial y}{\partial x} = 1 - y^2$ 이다. 이때 $y = \tanh(x)$ 의 값과 그 미분 값을 각각 그래프로 그리면 다음 그림처럼 된다.



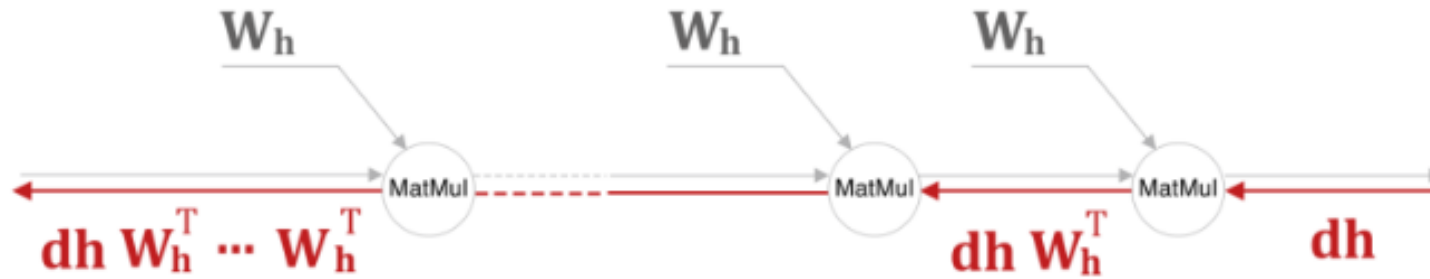
dtanh에 주목하면 그 값은 1.0이하이고, x가 0으로부터 멀어질수록 작아진다. 달리 말하면 이는 역전파에서는 기울기가 tanh 노드를 지날 때마다 값은 계속 작아진다는 뜻이다. 그래서 tanh 함수를 T번 통과하면 기울기도 T번 반복해서 작아지게 된다.

NOTE. RNN 계층의 활성화 함수로는 주로 tanh 함수를 사용하는데, 이를 ReLU로 바꾸면 기울기 소실을 줄일 수 있다.
「Improving performance of recurrent neural network with rely nonlinearity」 논문에서는 ReLU를 사용해 성능을 개선했다.

RNN의 문제점

▪ 기울기 소실과 기울기 폭발의 원인 - Matmul

다음으로 Matmul 노드에 주목해보자. 여기서는 이야기를 단순하게 하기 위해 tanh 노드를 무시하기로 한다.



다음처럼 상류로부터 dh 라는 기울기가 흘러들어왔을 때 이때 Matmul 노드에서의 역전파는 dhW_h^T 라는 행렬 곱으로 기울기를 계산한다. 그리고 같은 계산을 시계열 데이터의 시간 크기만큼 반복한다. 여기에서 주목할 점은 이 행렬 곱셈에서는 매번 똑같은 가중치인 W_h 가 사용된다는 점이다. 간단한 코드로 한 번 실험해보자.

RNN의 문제점

▪ Exploding gradient

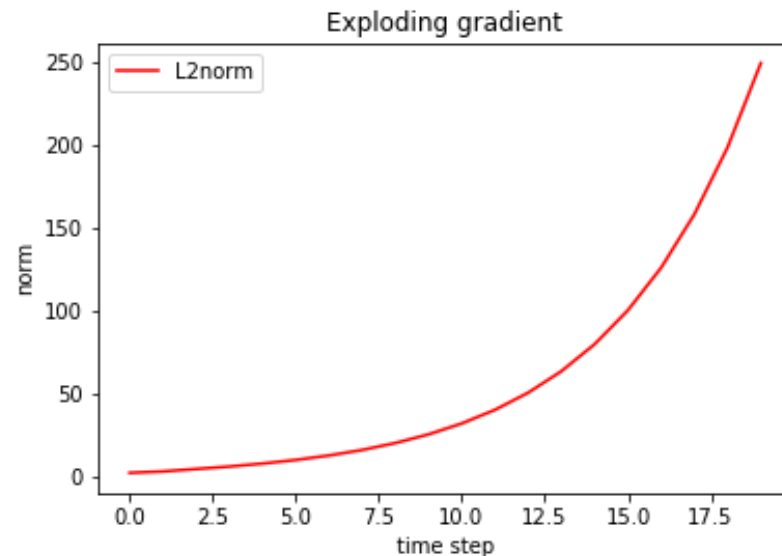
```
import numpy as np
import matplotlib.pyplot as plt

N = 2
H = 3
T = 20

dh = np.ones((N, H))
np.random.seed(3)
Wh = np.random.randn(H, H)

norm_list = list()
for t in range(T):
    dh = np.matmul(dh, Wh.T)
    norm = np.sqrt(np.sum(dh**2)) / N
    norm_list.append(norm)

plt.figure()
plt.plot(norm_list, label='L2norm', color = 'red')
plt.legend()
plt.title('Exploding gradient')
plt.xlabel('time step')
plt.ylabel('norm')
plt.savefig('gradient_explosion.png')
```



역전파의 Matmul 노드 수 (T)만큼 dh를 갱신했을 때, 기울기의 크기는 시간에 비례해 지수적으로 증가함을 알 수 있다. 이것이 바로 기울기 폭발 **exploding gradient**이다. 이러한 기울기 폭발이 일어나면 결국 Overflow를 일으켜 NaN(Not a Number) 같은 값을 발생시킨다.

RNN의 문제점

▪ Vanishing gradient

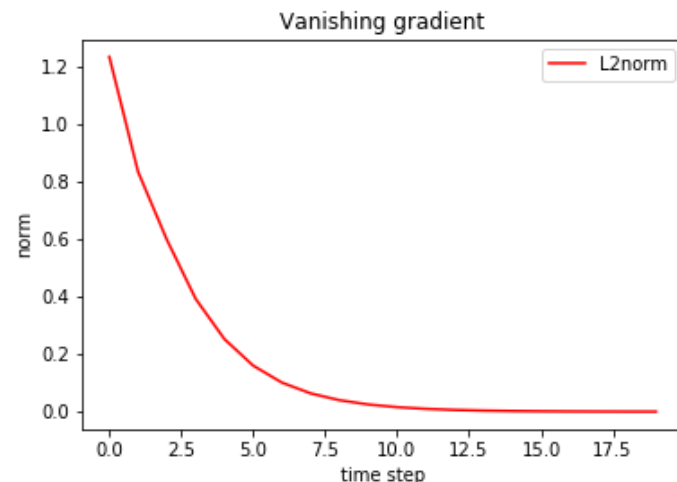
```
import numpy as np
import matplotlib.pyplot as plt

N = 2
H = 3
T = 20

dh = np.ones((N, H))
np.random.seed(3)
Wh = np.random.randn(H, H)

norm_list = list()
for t in range(T):
    dh = np.matmul(dh, Wh.T) / 2
    norm = np.sqrt(np.sum(dh**2)) / N
    norm_list.append(norm)

plt.figure()
plt.plot(norm_list, label='L2norm', color = 'red')
plt.legend()
plt.title('Vanishing gradient')
plt.xlabel('time step')
plt.ylabel('norm')
plt.show()
```

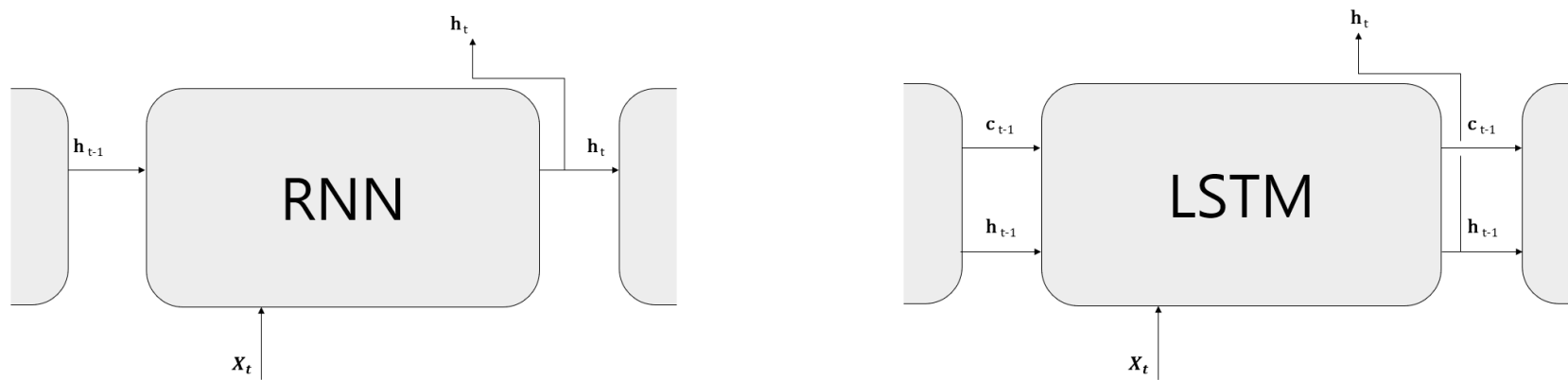


이번에는 초기값을 절반으로 줄이고 똑같이 돌려봤을 때, 기울기가 지수적으로 감소하는 것을 확인할 수 있다. 이것이 기울기 소실 **Vanishing gradient**이다. 이처럼 기울기 소실이 일어나서 기울기가 일정 수준 이하로 작아지면 가중치 매개변수가 더 이상 갱신되지 않는다.

LSTM

기울기 소실과 LSTM

이제 이러한 기울기 소실을 일으키지 않는다는 (혹은 일으키기 어렵게 한다는) LSTM 구조에 대해 살펴보고 이 구조를 개량한 GRU의 구조까지 살펴보자.



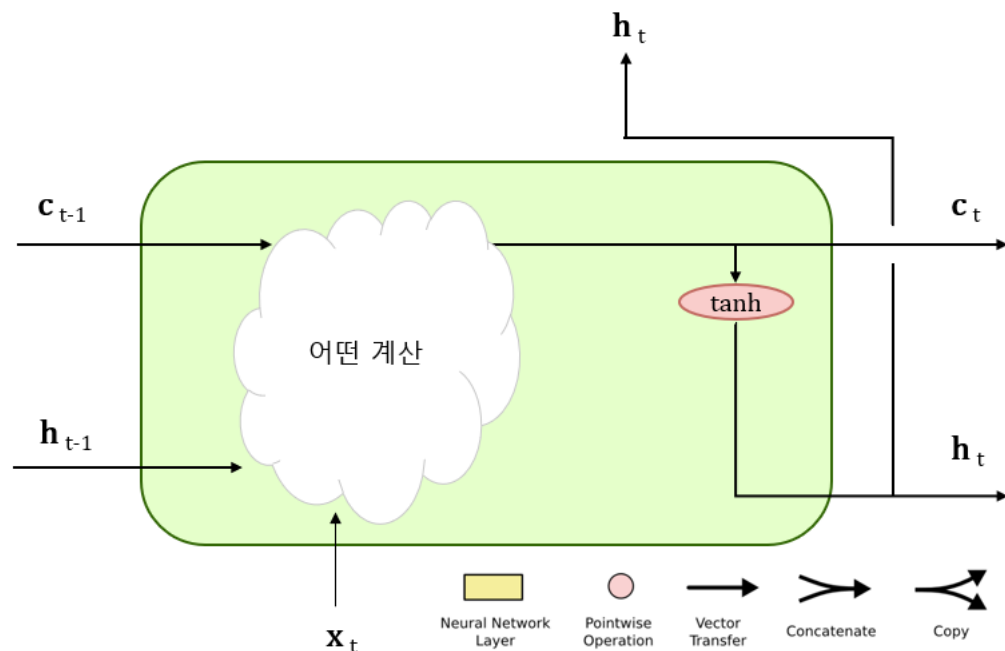
RNN 계층과 LSTM 계층 비교

다음은 RNN과 LSTM의 인터페이스를 비교한 그림이다. 그림에서 보듯 LSTM 계층의 인터페이스에는 **c**라는 경로가 있다는 차이가 있다. **c**를 기억 셀 **memory cell**(혹은 단순히 '셀')이라 하며, LSTM 전용의 **기억 매커니즘**이다. 기억 셀의 특징은 데이터를 자기 자신으로만 (LSTM 계층 내에서만) 주고받는다. 즉, LSTM 계층 내에서만 완결되고, 다른 계층으로는 출력하지 않습니다. 반면, LSTM의 은닉 상태 **h**는 RNN 계층과 마찬가지로 다른 계층으로 출력된다.

NOTE_ LSTM의 출력을 받는 쪽에서 보면 LSTM의 출력은 은닉 상태 벡터 **h**뿐이다. 그러므로 **c**의 존재 자체를 생각할 필요가 없다.

LSTM

▪ LSTM 계층 조립하기



기억 셀 c_t 를 바탕으로 은닉 상태 h_t 를 계산하는 LSTM 계층

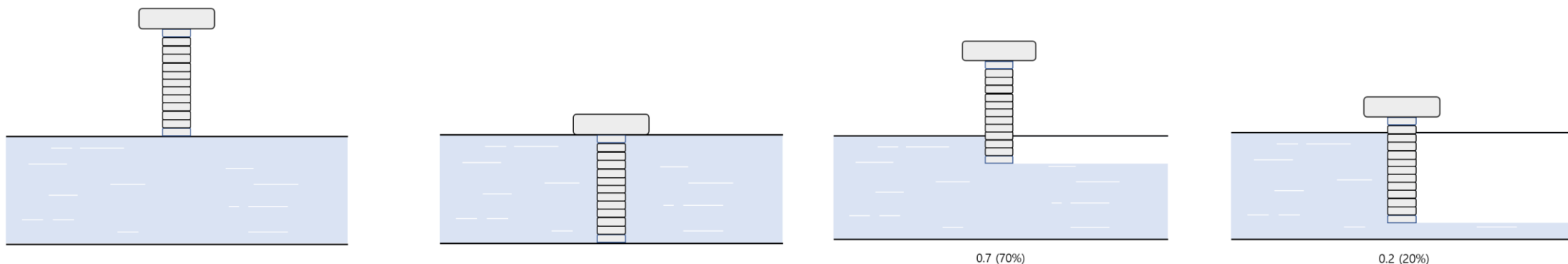
이제 LSTM의 구조를 차분히 알아보자. 앞서 이야기한 것처럼, LSTM에는 기억 셀 c_t 가 있다. 이 c_t 에는 시각 t 에서의 LSTM의 기억이 저장되어 있는데, 과거로부터 시각 t 까지에 필요한 모든 정보가 저장돼 있다고 가정하자. 그리고 필요한 정보를 모두 간직한 이 기억을 바탕으로 외부 계층에 은닉 상태 h_t 를 출력한다. 이때 출력하는 h_t 는 다음 그림과 같이 기억 셀의 값을 **tanh** 함수로 변환한 값이다. 그림처럼 현재의 기억 셀 c_t 는 3개의 입력 (c_{t-1} , h_{t-1} , x_t)으로부터 '어떤 계산'을 수행하여 구할 수 있다. 여기서 핵심은 갱신된 c_t 를 사용해 은닉상태 h_t 를 계산한다는 것이다. 또한 이 계산은 $h_t = \tanh(c_t)$ 인데, 이는 c_t 의 각 요소에 tanh 함수를 적용한다는 뜻이다. LSTM 구조에서의 핵심은 h_t 는 단기상태(short term state), c_t 는 장기상태(long term state)라고 볼 수 있다.

LSTM

▪ LSTM의 게이트^{gate}

진도를 더 나가기 전에, 이쯤에서 '게이트'라는 기능에 대해 얘기해보자. 게이트는 데이터의 흐름을 제어한다. 마치 다음 그림처럼 물의 흐름을 멈추거나 배출하는 것이 게이트의 역할이다.

그림1 비유하자면 게이트는 물의 흐름을 제어한다.



LSTM에서 사용하는 게이트는 '열기/닫기' 뿐 아니라 어느 정도 열지를 조절할 수 있다. 다음 그림처럼 게이트의 열림 상태는 0.0~1.0 사이의 실수로 나타난다. 그리고 그 값이 흐르는 물의 양을 결정한다. 여기서 중요한 것은 '게이트를 얼마나 열까'라는 것도 데이터로부터 자동으로 학습한다는 점이다.

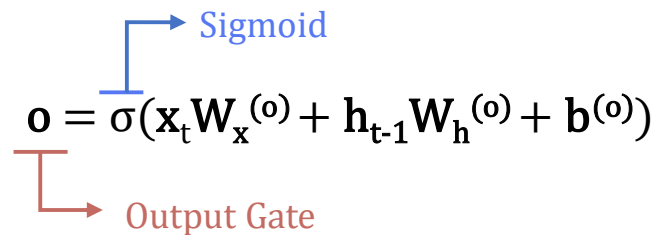
NOTE_ 게이트는 게이트의 열림 상태를 제어하기 위해서 전용 가중치 매개변수를 이용하며, 이 가중치 매개변수는 학습 데이터로부터 갱신된다. 참고로 게이트의 열림 상태를 구할 때는 시그모이드 함수를 사용하는데, 그 이유는 시그모이드 함수의 출력이 0.0~1.0 사이의 실수이기 때문이다.

LSTM

▪ output 게이트 - (1)

다시 LSTM 이야기로 돌아와보자. 바로 앞에서 은닉상태 \mathbf{h}_t 는 기억 셀 \mathbf{c}_t 에 단순히 \tanh 함수를 적용했을 뿐이라고 설명했다. 이번에는 $\tanh(\mathbf{c}_t)$ 에 게이트를 적용하는 걸 생각해보자. 즉, $\tanh(\mathbf{c}_t)$ 의 각 원소에 대해 '그것이 다음 시각의 은닉 상태에 얼마나 중요한가'를 조정한다.

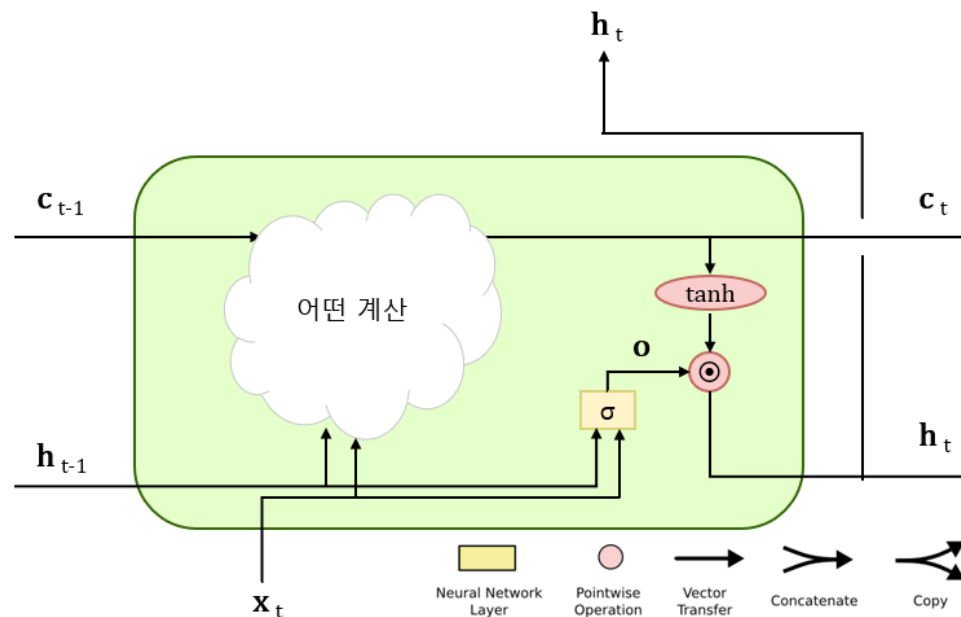
output 게이트의 열림 상태 (다음 몇 %만 흘려보낼까)는 입력 \mathbf{x}_t 와 이전 상태 \mathbf{h}_{t-1} 로부터 구한다. 이때의 계산은 다음과 같다.

$$\mathbf{o} = \sigma(\mathbf{x}_t \mathbf{W}_x^{(o)} + \mathbf{h}_{t-1} \mathbf{W}_h^{(o)} + \mathbf{b}^{(o)})$$


위의 식을 보게되면 RNN 계층의 계산에서 \tanh 가 아닌 Sigmoid를 사용했다는 점만이 다르다는 것을 확인할 수 있다.

LSTM

▪ output 게이트 - (2)



Output 게이트가 추가된 LSTM 계산 그래프

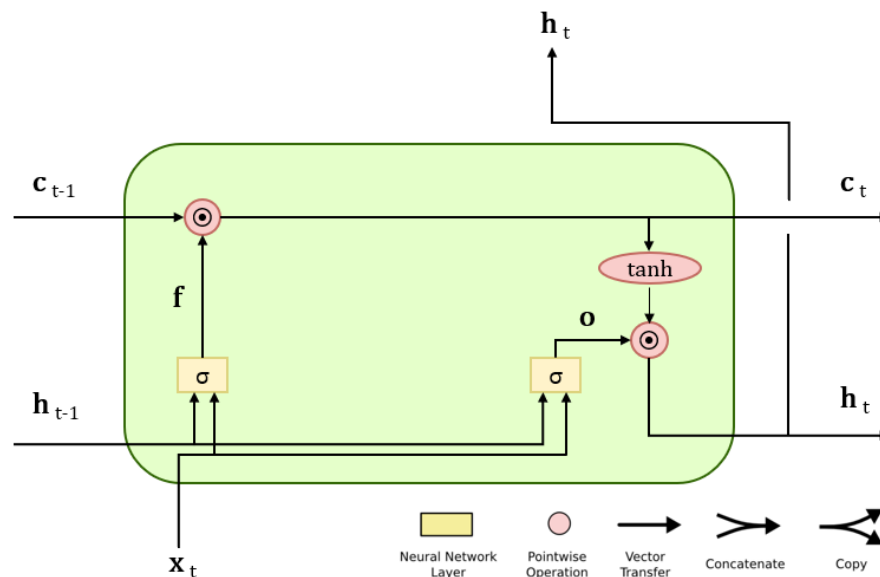
output 게이트에서 수행하는 식을 σ 로 표기했다. 그리고 σ 의 출력을 o 라고 하면 h_t 는 o 와 $\tanh(c_t)$ 의 곱으로 계산된다. 이때 말하는 '곱'이란 원소별 **element-wise**이며, 이것을 아다마르 곱 **Hardward product**이라고도 한다. 아다마르 곱을 기호로는 \odot 로 나타내며 다음과 같은 계산을 수행한다.

$$h_t = o \odot \tanh(c_t)$$

NOTE_ \tanh 의 출력은 -1.0~1.0의 실수고, 이 -1.0~1.0의 수치를 인코딩된 '정보'의 강약(정도)을 표시한다고 해석할 수 있다. 한편 시그모이드 함수의 출력은 0.0~1.0의 실수이며, 데이터를 얼마만큼 통과시킬지를 정하는 비율이 된다. 주로 게이트에서는 시그모이드 함수가, 실질적인 '정보'를 지니는 데이터에는 \tanh 함수가 활성화 함수로 사용된다.

LSTM

- **forget 게이트**



forget 게이트까지 추가된 LSTM 계산 그래프

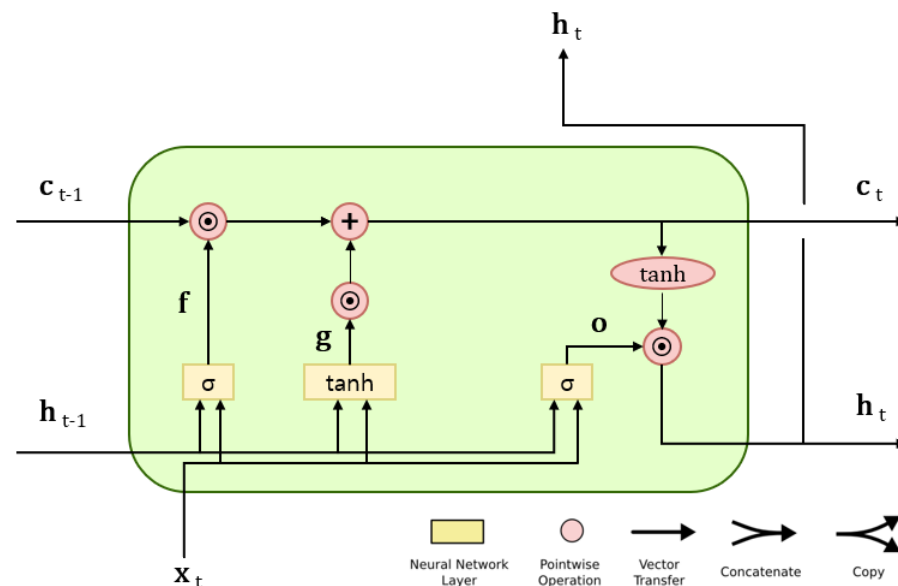
다음으로 할 일은 기억 셀에 '무엇을 잊을까'를 지시하는 것이다. 이것도 역시 게이트를 사용해 해결한다. c_{t-1} 의 기억 중에서 불필요한 기억을 잊게 해주는 게이트를 **forget 게이트**(망각 게이트)라고 한다. Forget 게이트가 수행하는 일련의 계산을 σ 노드로 표기했다. σ 안에서는 다음 식의 계산을 수행한다.

$$\mathbf{f} = \sigma(\mathbf{x}_t \mathbf{W}_x^{(f)} + \mathbf{h}_{t-1} \mathbf{W}_h^{(f)} + \mathbf{b}^{(f)})$$

그리고 이 \mathbf{f} 와 이전 기억 셀인 \mathbf{c}_{t-1} 과의 아다마르 곱, 즉 $\mathbf{c}_t = \mathbf{f} \odot \mathbf{c}_{t-1}$ 을 계산하여 \mathbf{c}_t 를 구한다.

LSTM

▪ 새로운 기억 셀



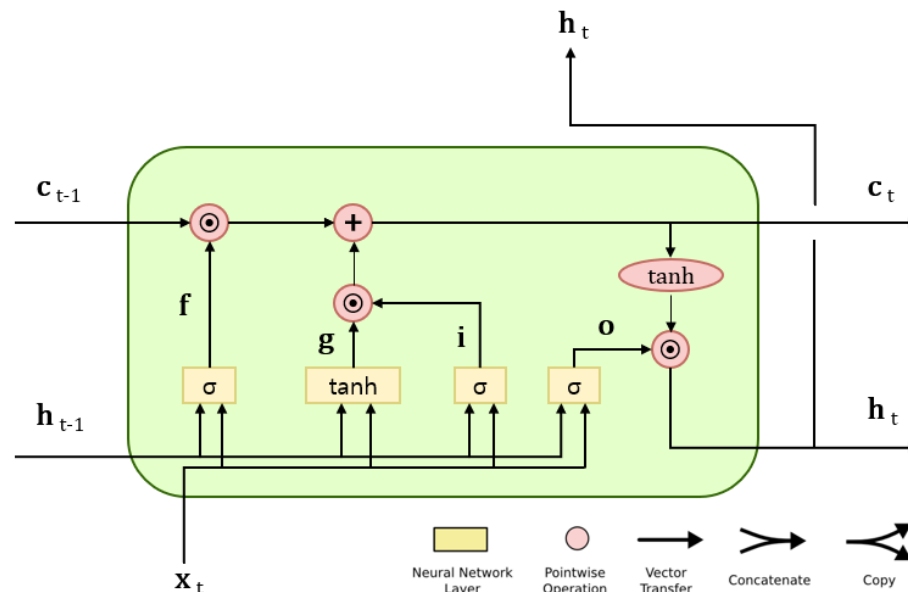
새로운 기억 셀에 필요한 정보가 추가된 LSTM 계산 그래프

forget 게이트를 거치면서 이전 시각의 기억 셀로부터 잊어야 할 기억이 삭제됐다. 이 상태로는 기억 셀이 잊는것 밖에 하지 못하므로 새로 기억해야 할 정보를 추가해야한다. 그러기 위해서 위의 그림과 같이 **tanh 노드를 추가**한다. ('정보'가 담겨있으므로) 위의 그림에서 보듯 tanh 노드가 계산한 결과가 이전 시각의 기억 셀 c_{t-1} 에 더해진다. 기억 셀이 새로운 '정보'가 추가된 것이다. 주의할 점은 이 tanh 노드는 '게이트'가 아니며, 새로운 '정보'를 기억 셀에 추가하는 것이 목적이다. 따라서 활성화 함수로는 시그모이드 함수가 아닌 tanh 함수가 사용된다. 이제 잊는 것 뿐만이 아닌, 새로운 정보까지 추가가 되었다.

$$g = \tanh(\mathbf{x}_t \mathbf{W}_x^{(g)} + \mathbf{h}_{t-1} \mathbf{W}_h^{(g)} + \mathbf{b}^{(g)})$$

LSTM

▪ input 게이트



input 게이트까지 추가된 LSTM 계산 그래프

마지막으로 새로운 정보가 들어있는 g 에 게이트를 하나 추가할 생각이다. 여기에서 새롭게 추가하는 게이트를 **input 게이트**라고 한다. Input 게이트를 g 의 각 원소가 새로 추가되는 정보로써의 가치가 얼마나 큰지를 판단한다. 즉, 새로운 정보를 무비판적으로 수용하는 것이 아니라, 적절히 취사선택하는 것이 이 게이트의 역할이다.

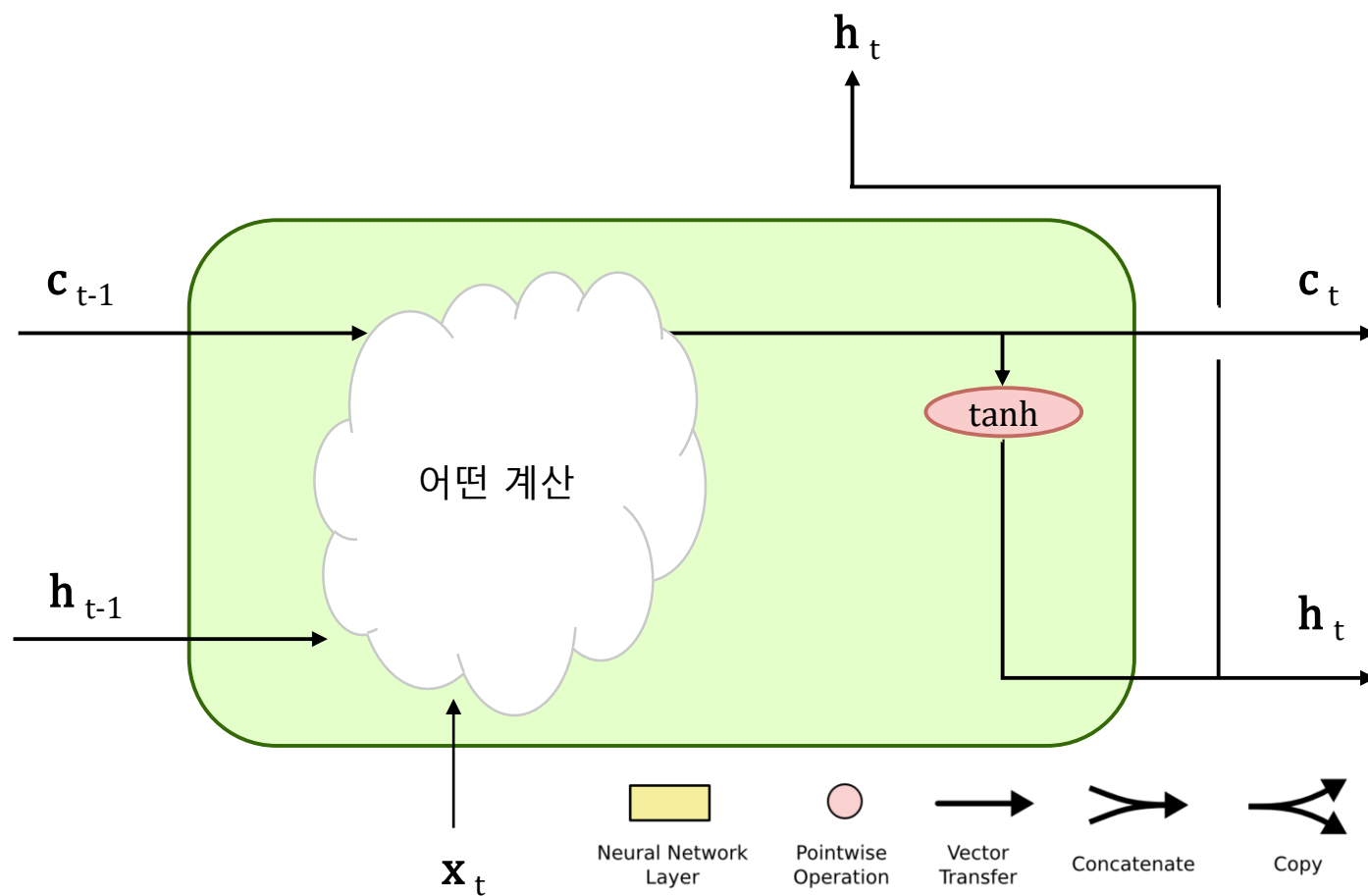
$$i = \sigma(x_t W_x^{(i)} + h_{t-1} W_h^{(i)} + b^{(i)})$$

그런 다음 i 와 g 의 아다마르 곱 결과를 기억 셀에 추가한다. 이상이 LSTM 안에서 이뤄지는 처리이다.

LSTM Review

Long Short Term Memory (LSTM)

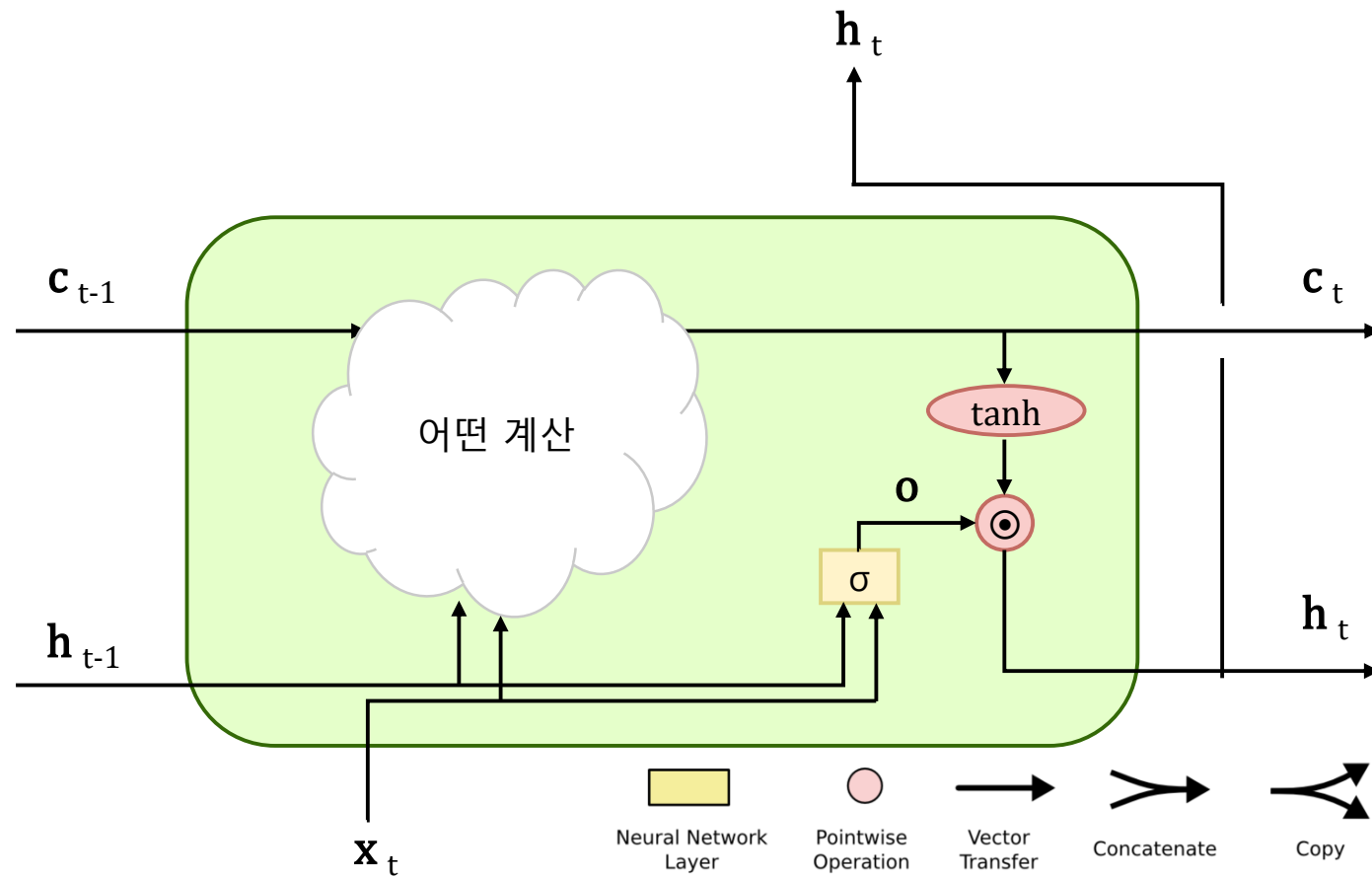
기억 셀 c_t 를 바탕으로 은닉상태 h_t 를 계산하는 LSTM 계층



LSTM Review

Long Short Term Memory (LSTM)

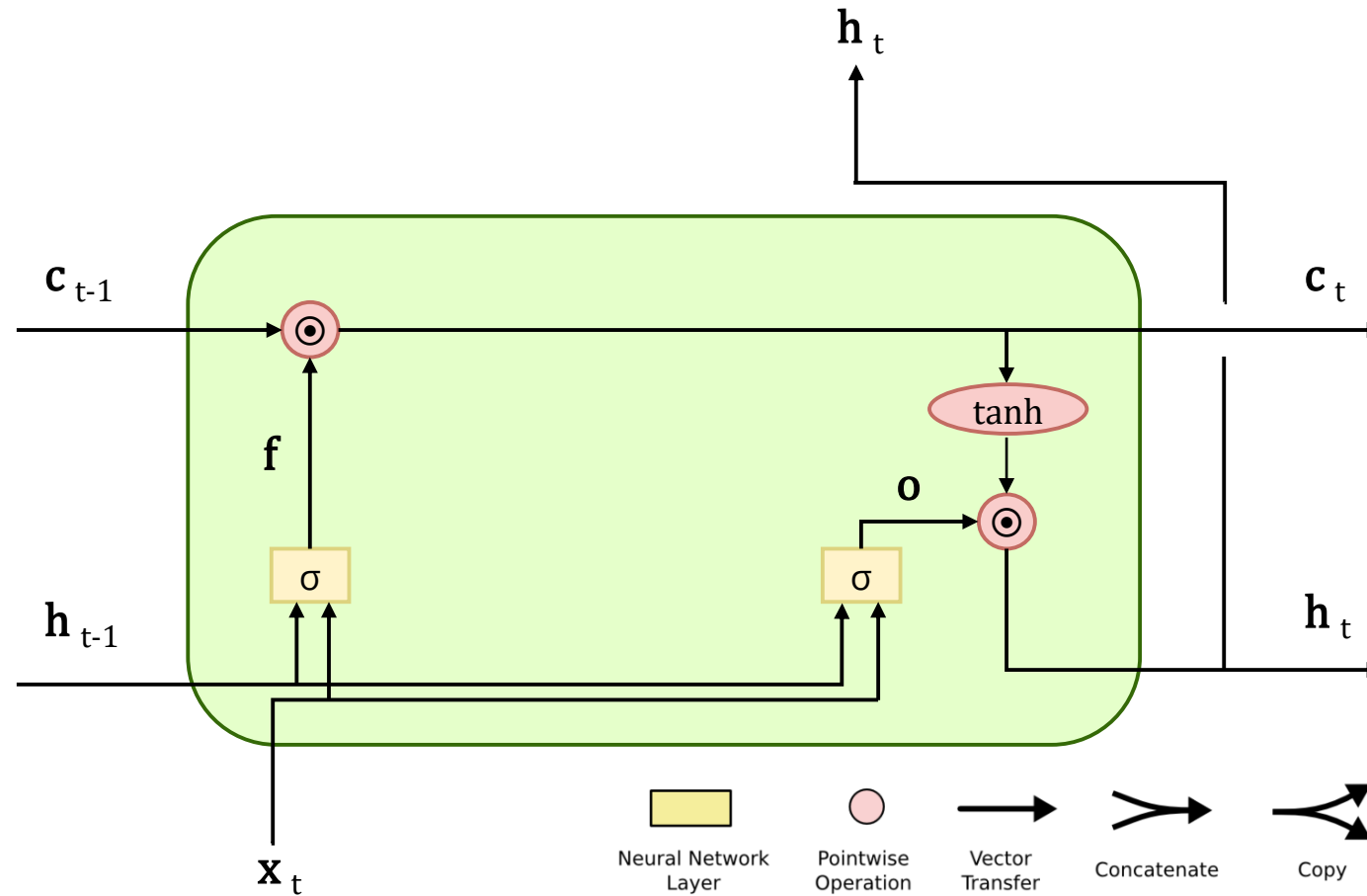
output 게이트 추가 (o gate)



LSTM Review

Long Short Term Memory (LSTM)

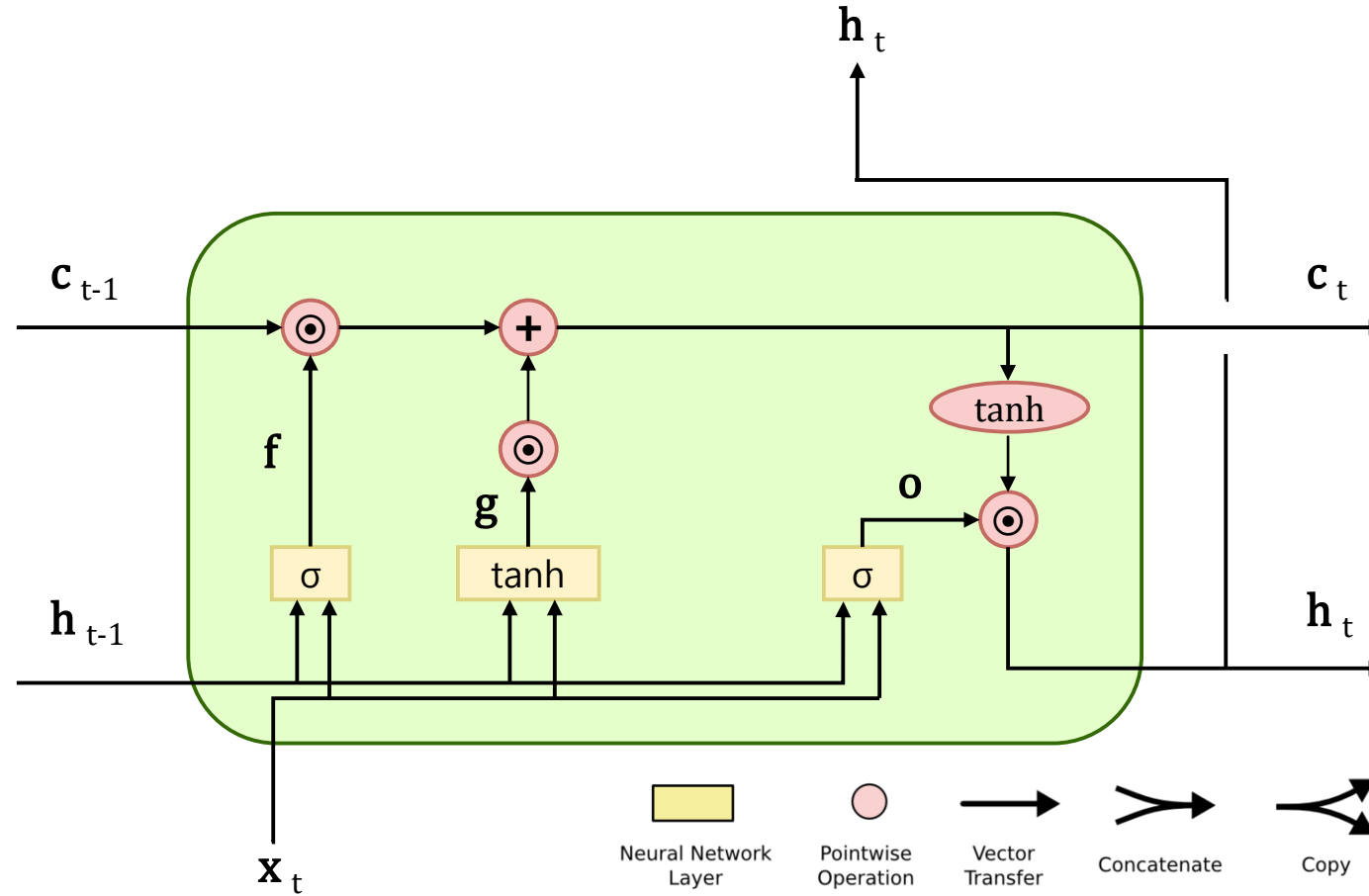
forget 게이트 추가 (f gate)



LSTM Review

Long Short Term Memory (LSTM)

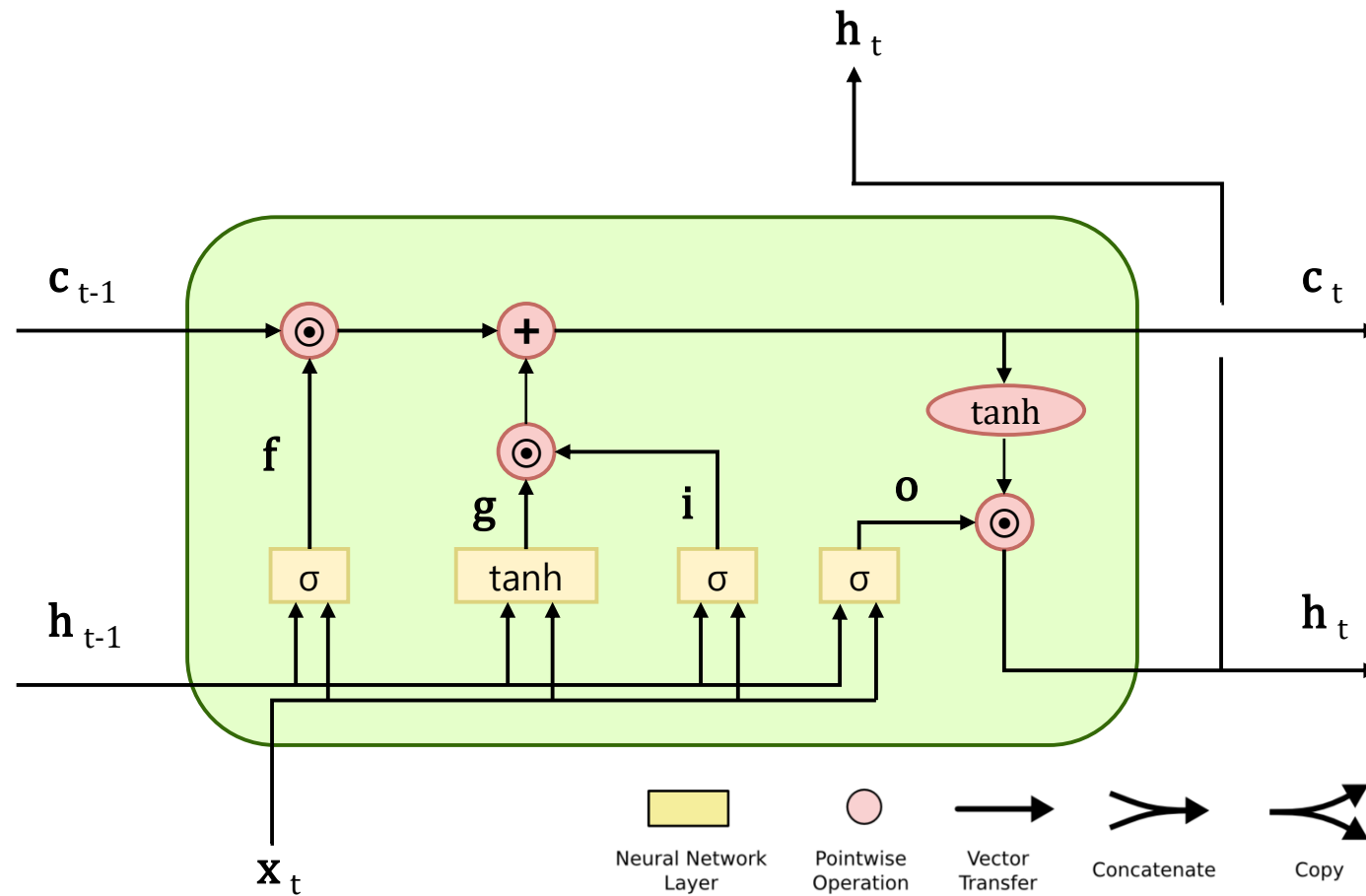
새로운 기억 셀에 필요한 정보를 추가 (g gate)



LSTM Review

Long Short Term Memory (LSTM)

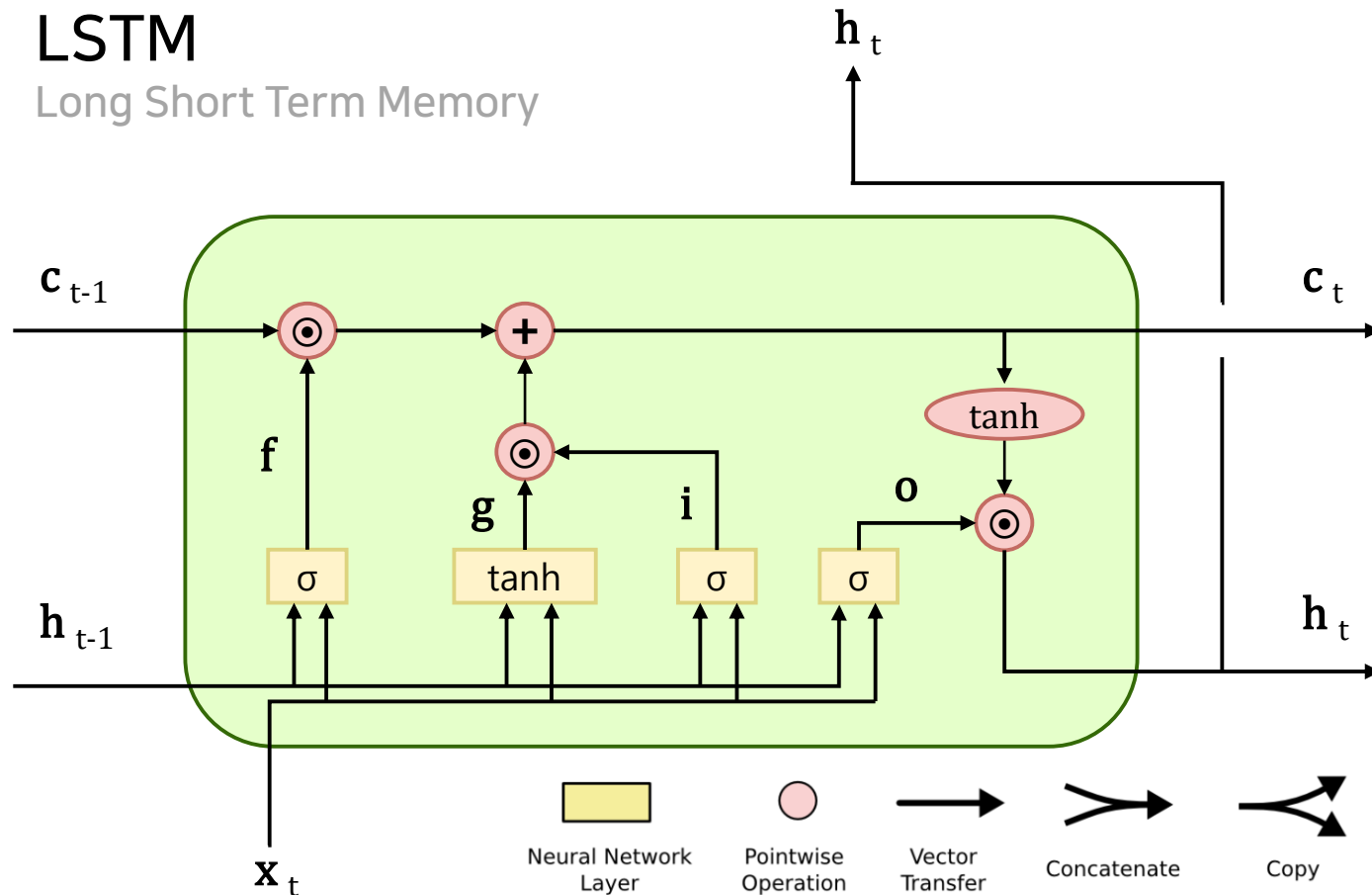
Input 게이트 추가 (i gate)



LSTM Review

Long Short Term Memory (LSTM)

LSTM의 계산 그래프



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

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

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

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

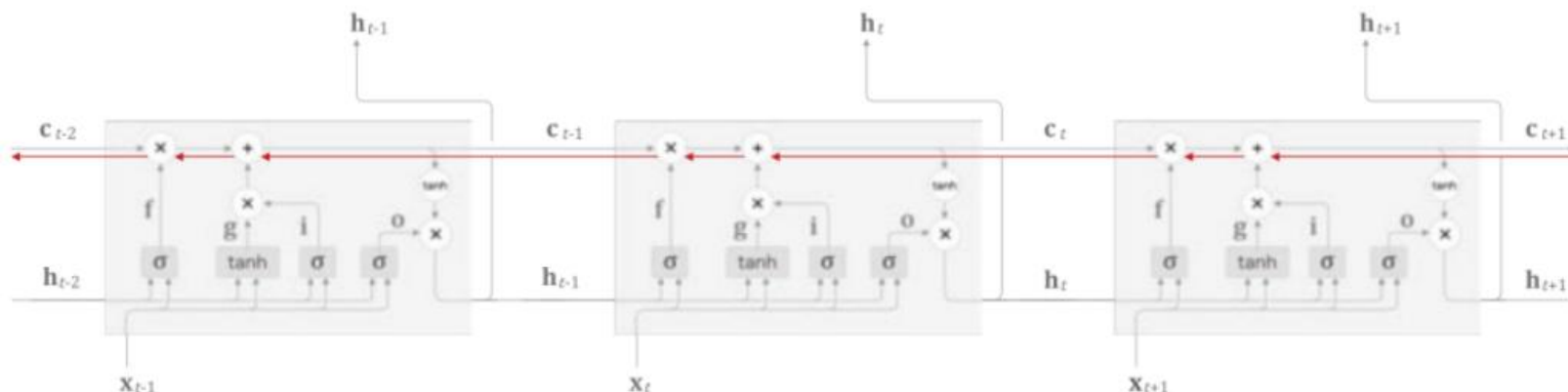
$$c_t = f \odot c_{t-1} + g \odot i$$

$$h_t = o \odot \tanh(c_t)$$

LSTM

▪ LSTM의 기울기 흐름

LSTM의 구조는 설명했지만, 이것이 어떤 원리로 기울기 소실을 없애주는 걸까? 그 원리는 기억 셀 c 의 역전파에 주목하면 볼 수 있다.



다음은 기억 셀에만 집중하여, 그 역전파의 흐름을 그린 것이다. 이때 기억 셀의 역전파에서는 '+'와 'X' 노드만을 지나게 된다. '+' 노드는 상류 기울기를 그대로 흘릴 뿐이므로 남는 것은 'X' 노드인데, 이 노드는 **'행렬 곱'이 아닌 아다마르 곱을 계산한다**. RNN의 역전파에서는 똑같은 가중치 행렬을 사용해서 '행렬 곱'을 반복했고, 그래서 기울기 소실 혹은 폭발이 일어났다. 반면 이번 LSTM의 역전파에서는 '행렬 곱'이 아닌 '원소별 곱'이 이뤄지고, 매 시각 다른 게이트 값을 이용해 원소별 곱을 계산한다. 이처럼 **매번 새로운 게이트 값**을 이용하므로 곱셈의 효과가 누적되지 않아 기울기 소실이 일어나기 어려운 것이다.

NOTE_ 위 그림의 'X' 노드의 계산은 forget 게이트가 제어한다. 역전파 계산시 forget 게이트의 출력과 상류 기울기의 곱이 계산되므로 forget 게이트가 '잊어야 한다'고 판단한 기억 셀의 원소에 대해서는 그 기울기가 작아지고, '잊어서는 안 된다'고 판단한 원소에 대해서는 그 기울기가 약화되지 않은 채로 과거 방향으로 전해진다. 따라서 중요한 정보의 기울기는 소실 없이 전파된다.

부록 - Hardmard Product

- Hardmard Product

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \circ \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix} = \begin{bmatrix} a_{11} b_{11} & a_{12} b_{12} & a_{13} b_{13} \\ a_{21} b_{21} & a_{22} b_{22} & a_{23} b_{23} \\ a_{31} b_{31} & a_{32} b_{32} & a_{33} b_{33} \end{bmatrix}$$

일반 행렬 곱은 $m \times n$ 행렬과 $n \times p$ 꼴의 행렬을 곱하지만, Hardmard product는 $m \times n$ 과 $m \times n$ 의 같은 꼴을 가지는 행렬끼리 같은 위치의 원소끼리 각각 곱한다