

Long Short-Term Memory network – from Scratch

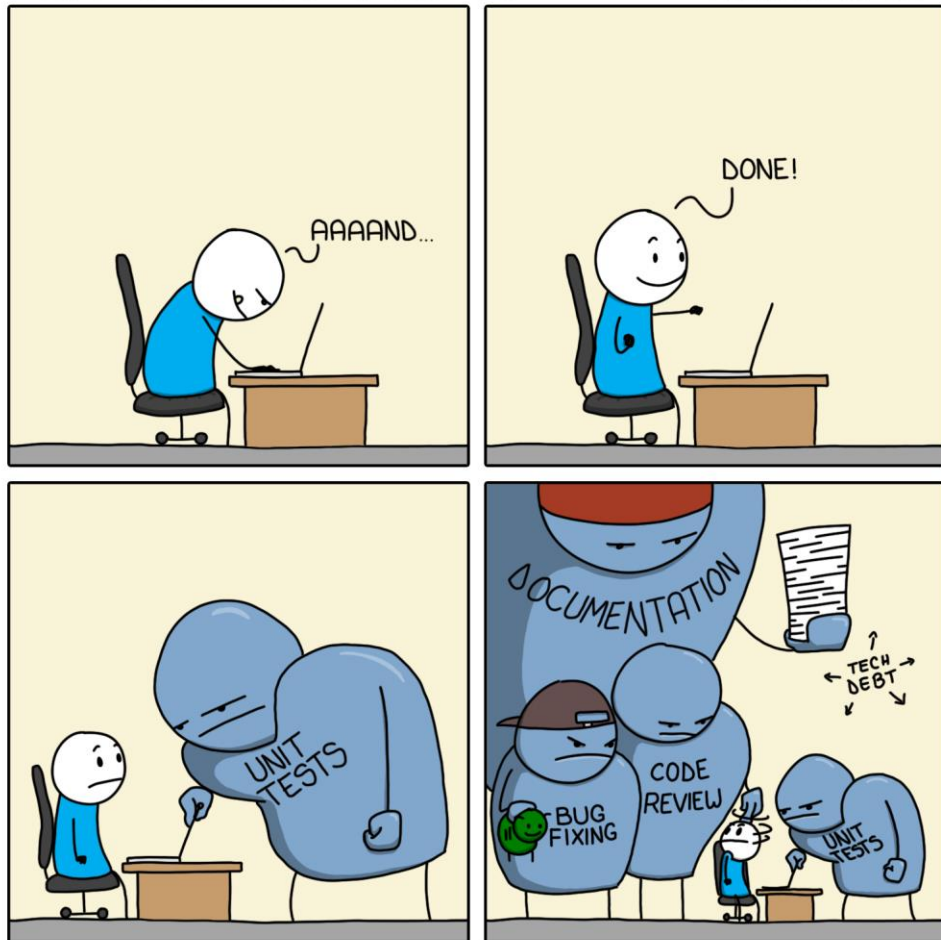


<https://www.analyticsvidhya.com>



FEATURE COMPLETE

MONKEYUSER.COM



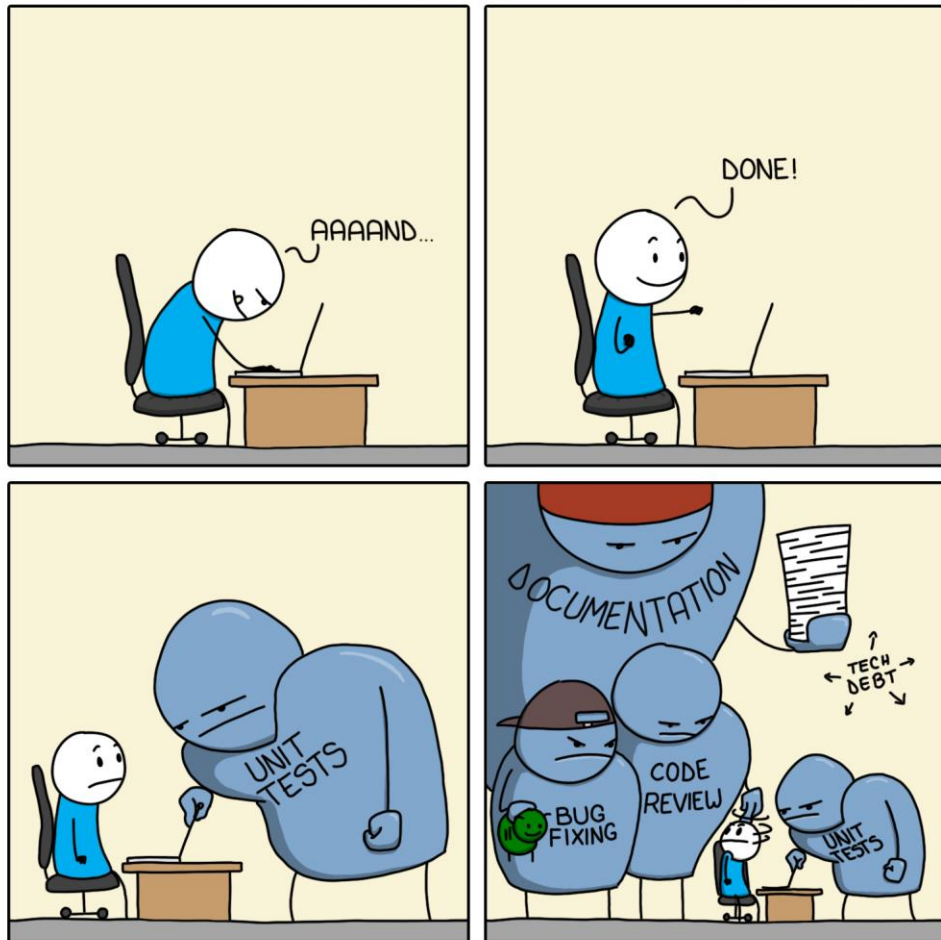
you should be fluent with:

- *basic OOP* (methods, classes, inheritance)
- *Linear algebra* (dot product, inner product, outer product)
- *derivatives* (gradient)
- +
- *Codes from “RNN from Scratch”*



FEATURE COMPLETE

MONKEYUSER.COM



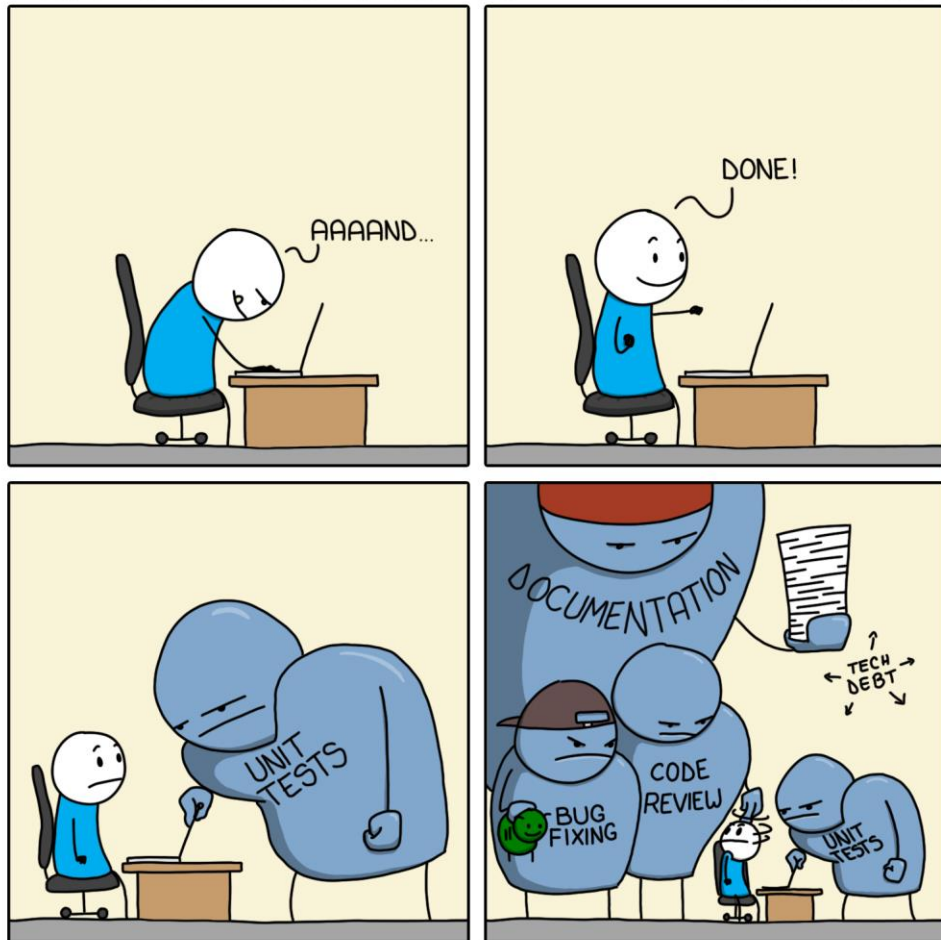
outline:

- the idea
- from a classical RNN cell to a LSTM
- **BackPropagation Through Time**
- full backpropagation
- modifying the SGD optimizer
- running and testing the package



FEATURE COMPLETE

MONKEYUSER.COM



outline:

- *the idea*

- *from a classical RNN cell to a LSTM*
- *BackPropagation Through Time*
- *full backpropagation*
- *modifying the SGD optimizer*
- *running and testing the package*



- time series analysis (prediction and forecasting)
- speech recognition
- anomaly detection

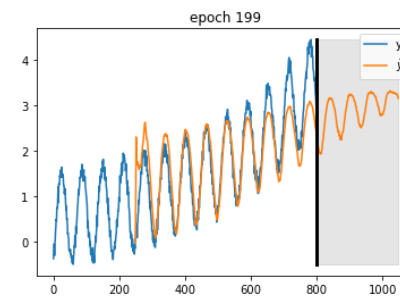
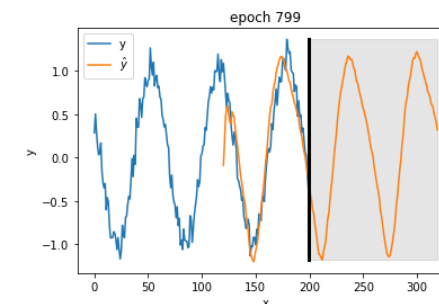
new:

- **long-term** and **short-term** memory
- dealing with vanishing/exploding gradient
- invented 1997 by Sepp Hochreiter und Jürgen Schmidhuber

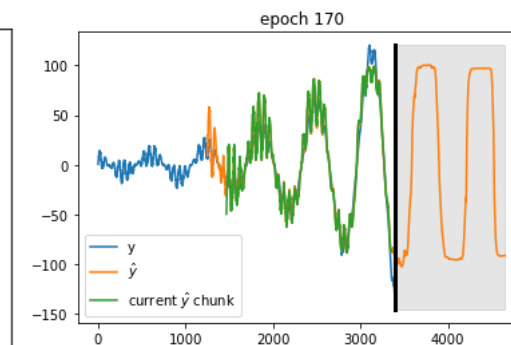
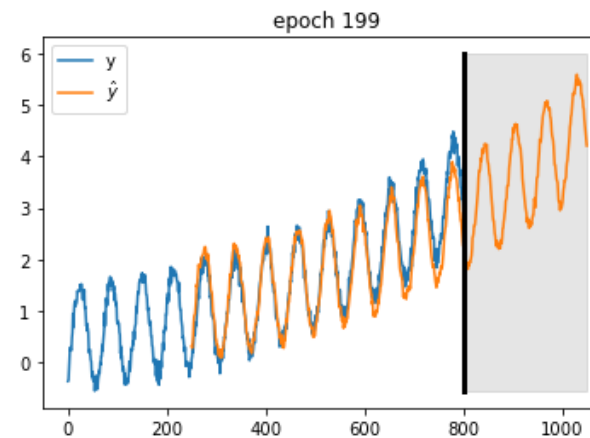
Aim of the lecture:

building a **many – to – many**,
one feature LSTM from our **previous RNN**

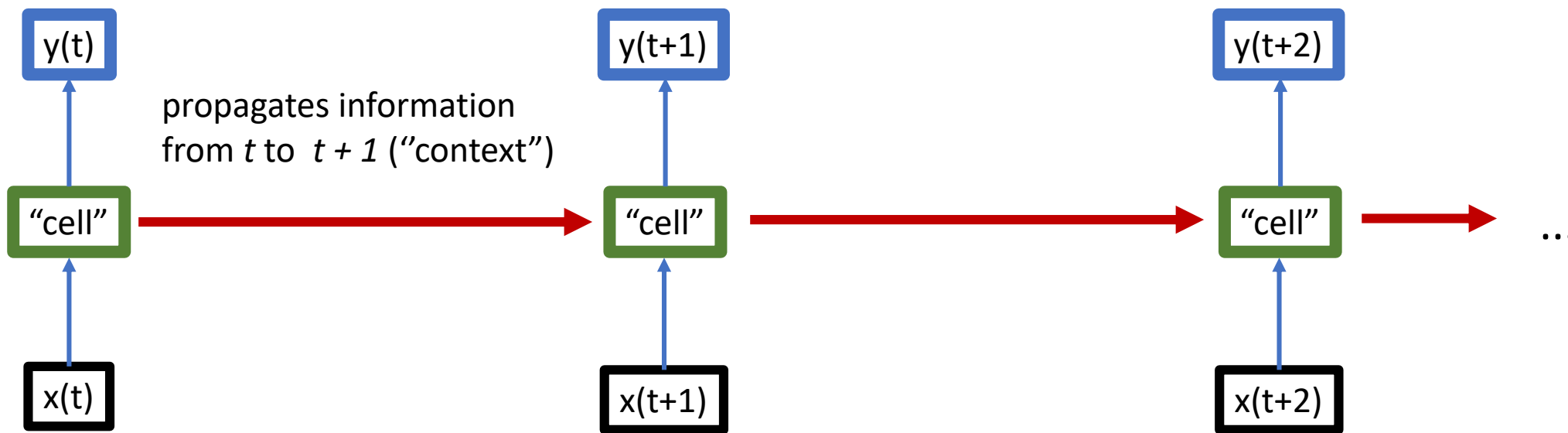
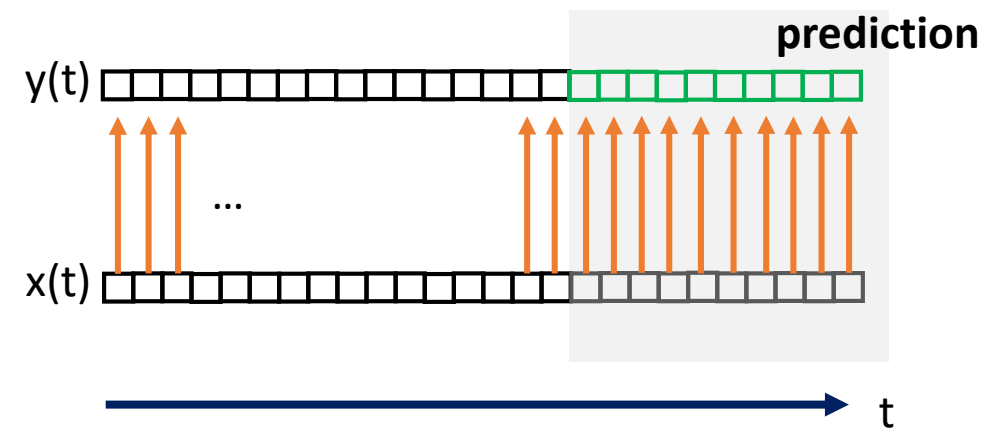
classical RNN



LSTM

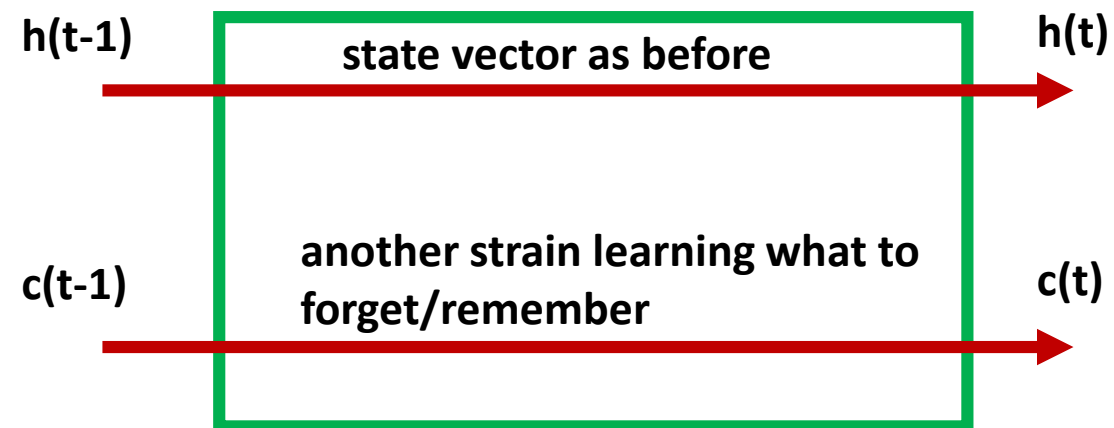
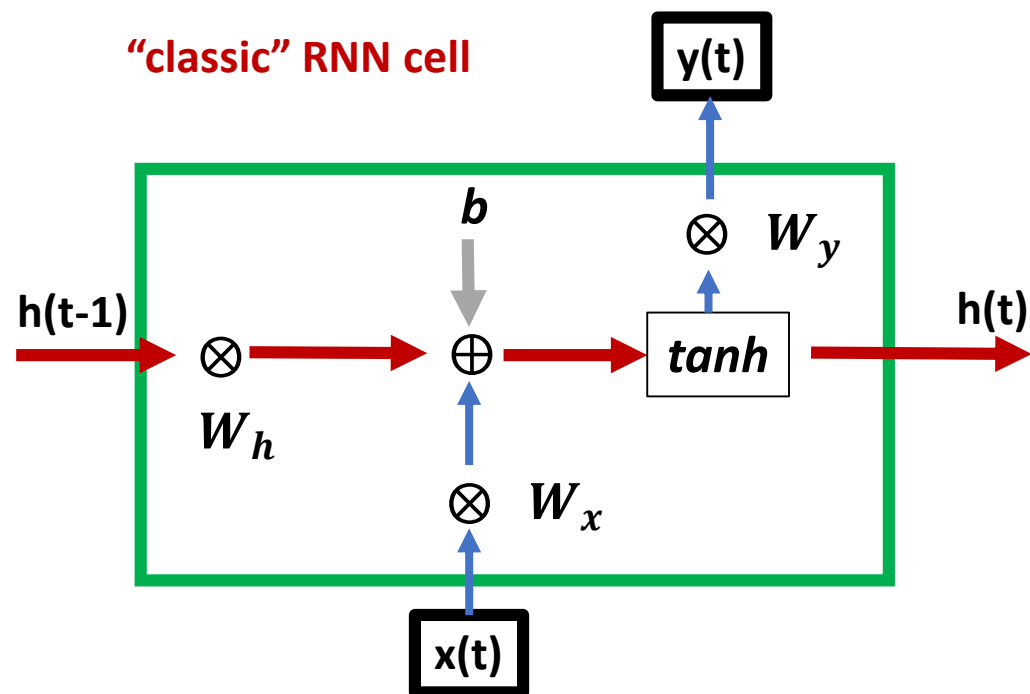


(adding more noise)



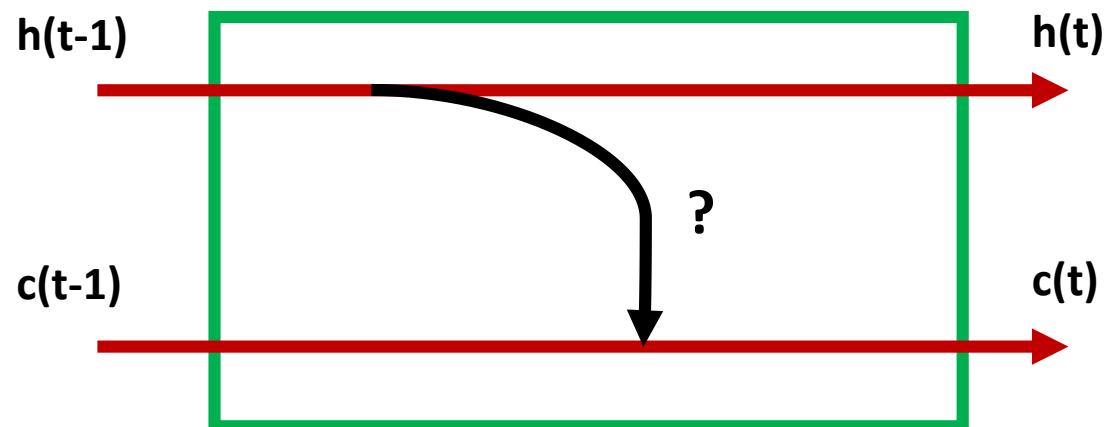
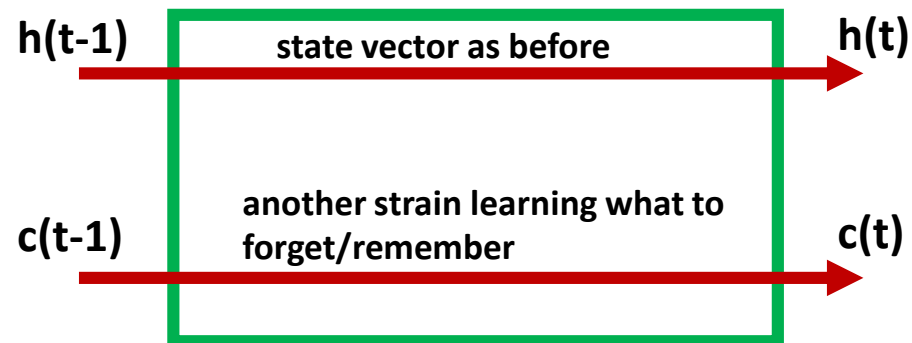
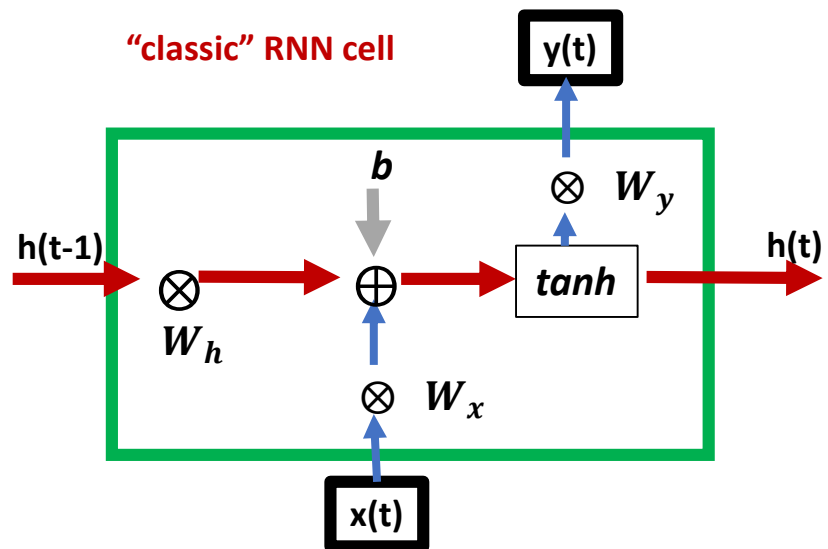


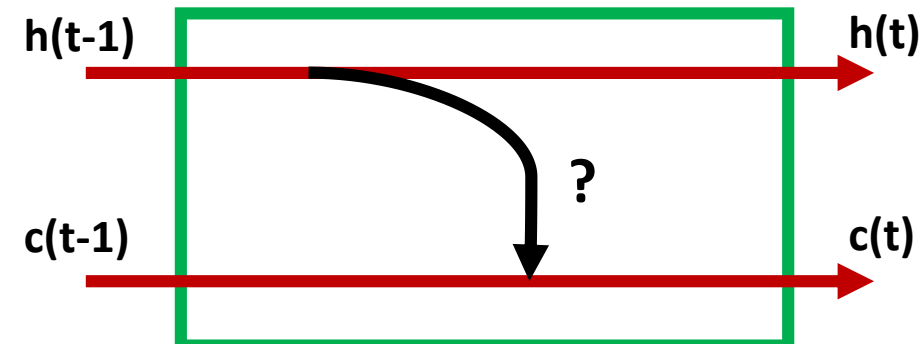
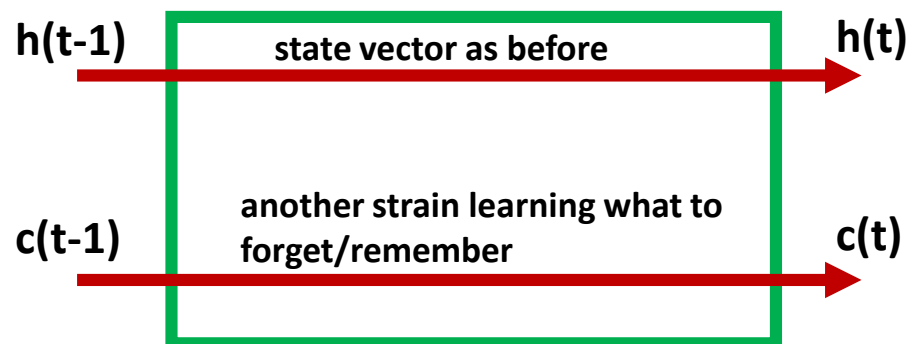
“classic” RNN cell



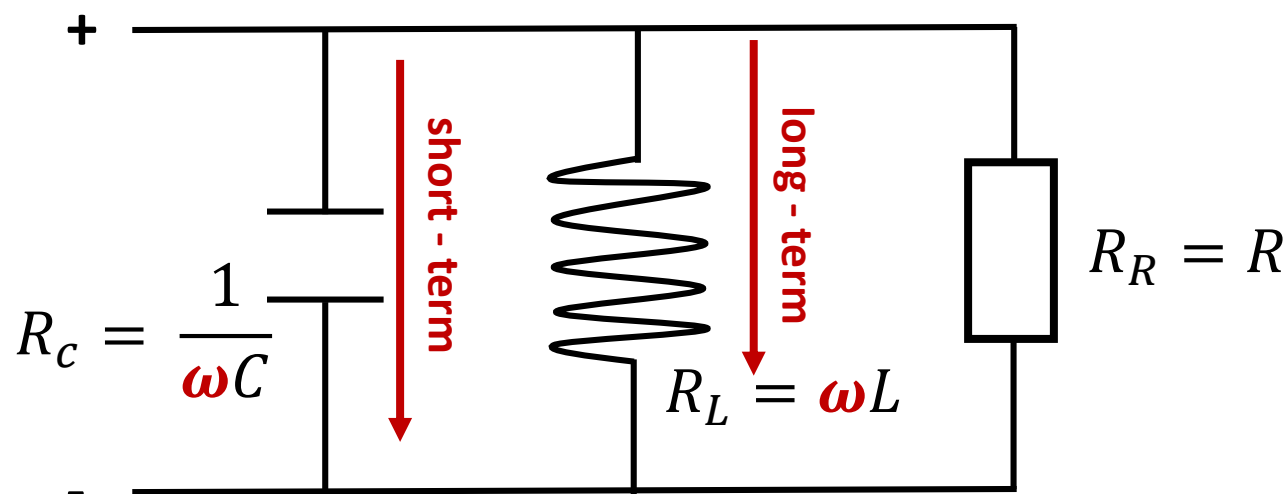


“classic” RNN cell





electrical circuits:

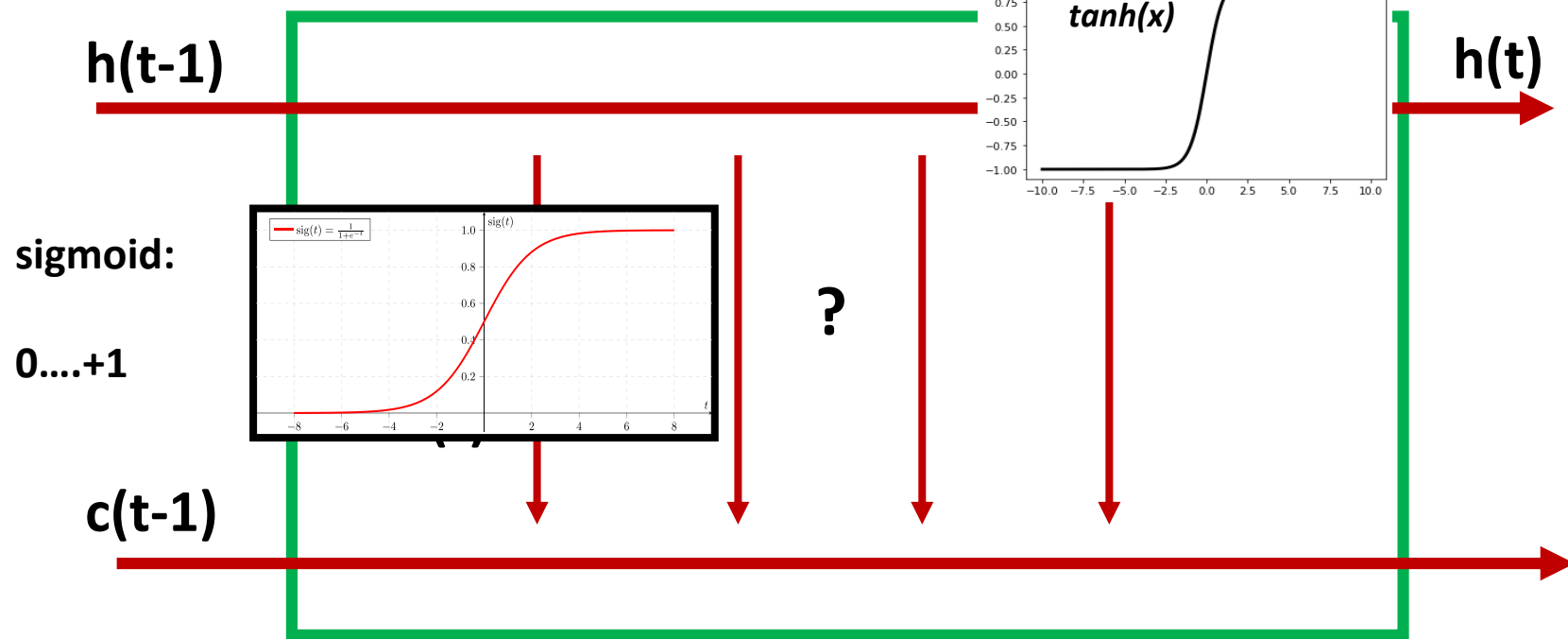
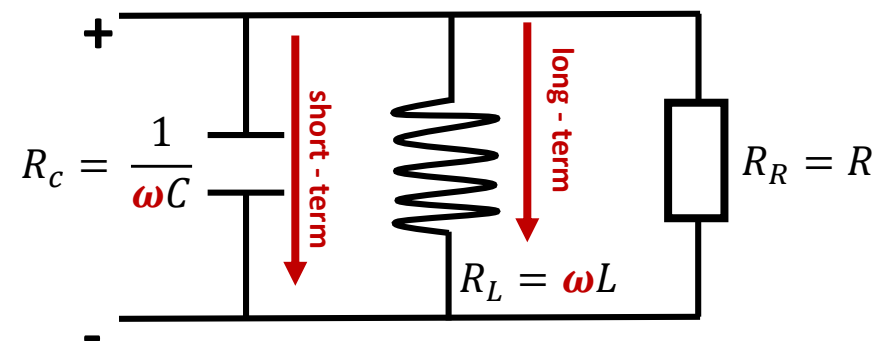
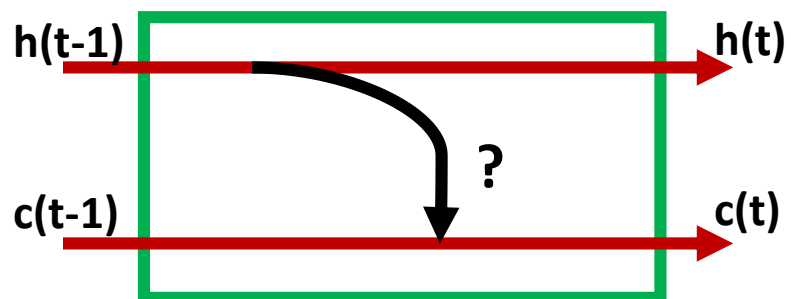


$$\text{AC: } I(t) = I_0 e^{i(\omega t + \varphi)}$$

R_c : passes **short** -term changes

R_L : passes **long** -term changes

$$\frac{1}{R_{tot}} = \frac{1}{R_R} + \frac{1}{R_c} + \frac{1}{R_L}$$



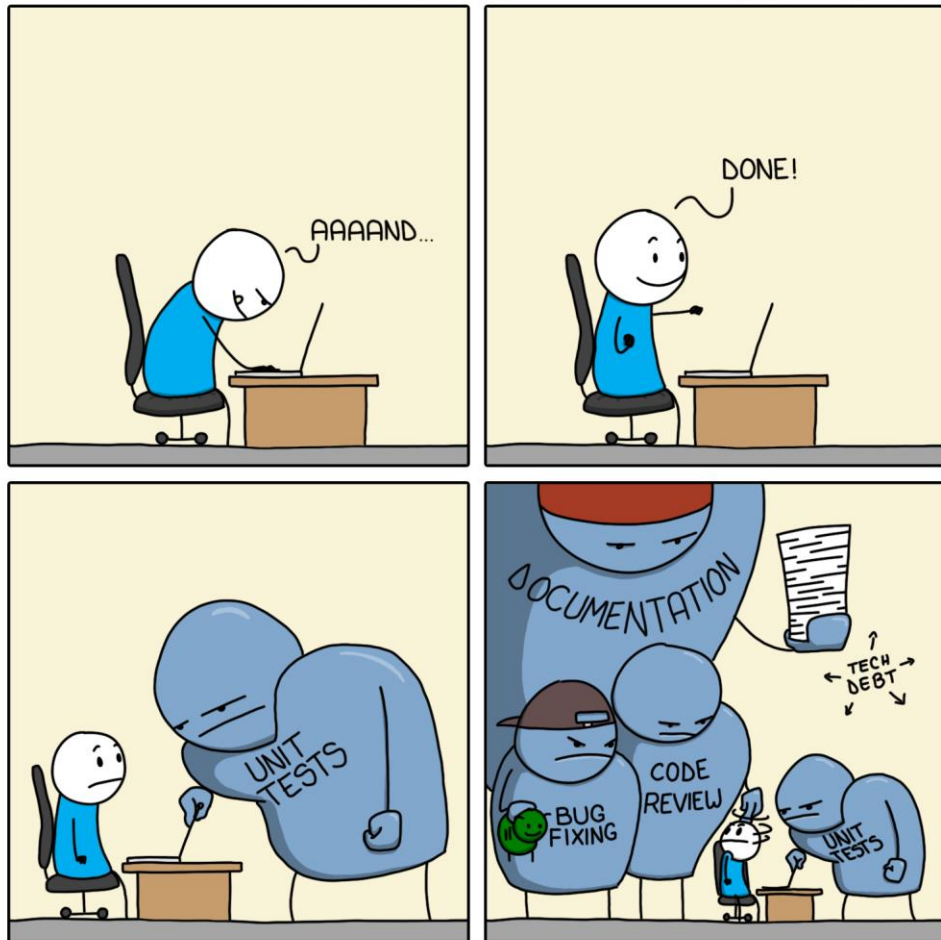
sigmoid:
learning whether
information has to
be passed on

tanh:
state vector as for RNN



FEATURE COMPLETE

MONKEYUSER.COM



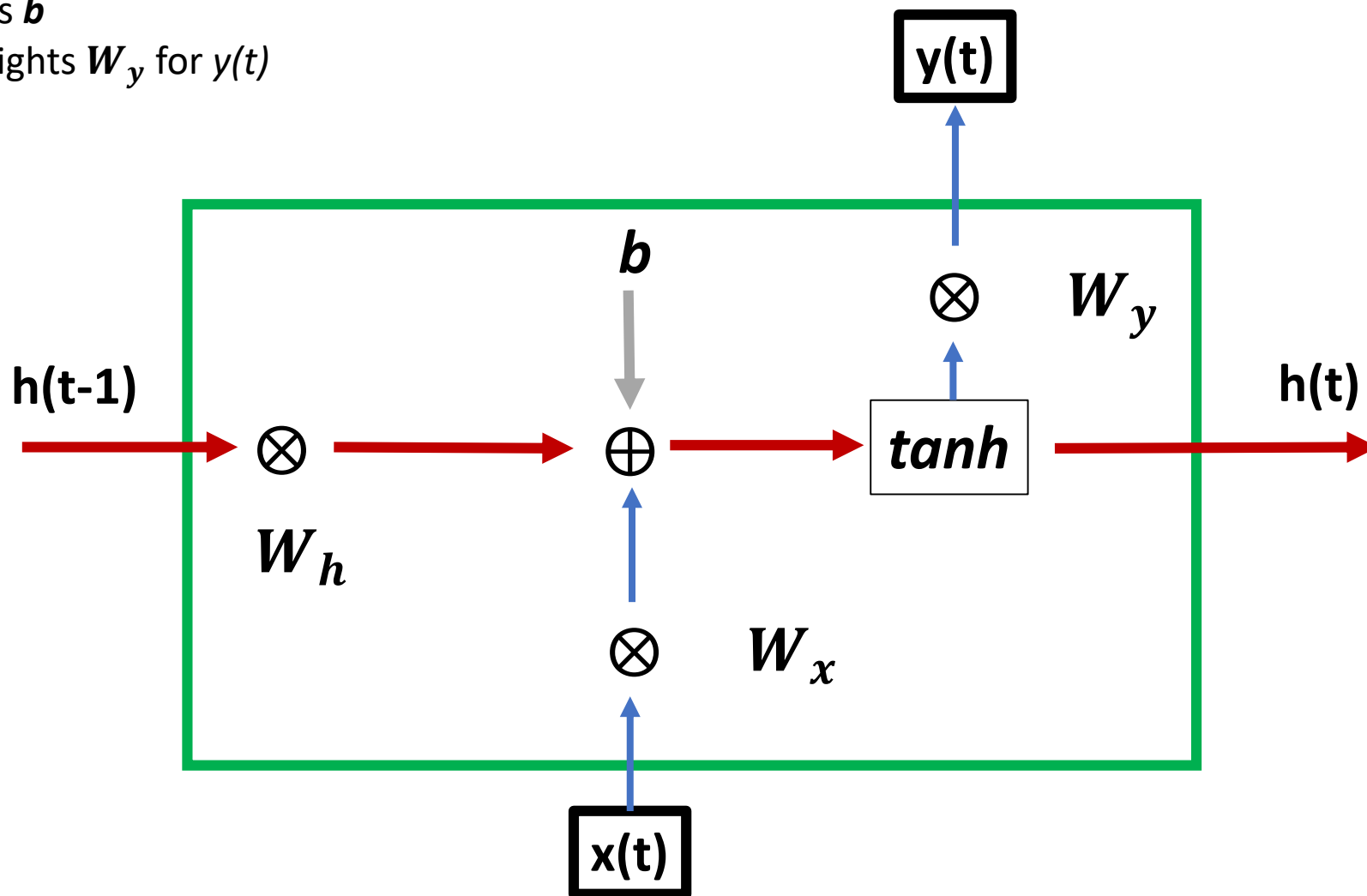
outline:

- *the idea*
- **from a classical RNN cell to a LSTM**
- *BackPropagation Through Time*
- *full backpropagation*
- *modifying the SGD optimizer*
- *running and testing the package*

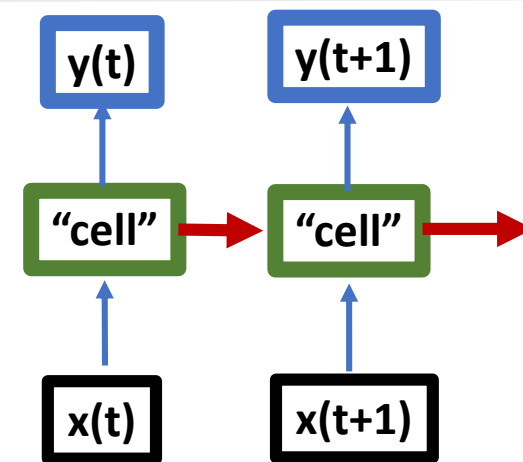


Long Short-Term Memory network - from Scratch

- **state vector** $h(t)$, i. e. the “memory” the system has along time
- weights W_x for $x(t)$
- weights W_h for $h(t)$
- bias b
- weights W_y for $y(t)$



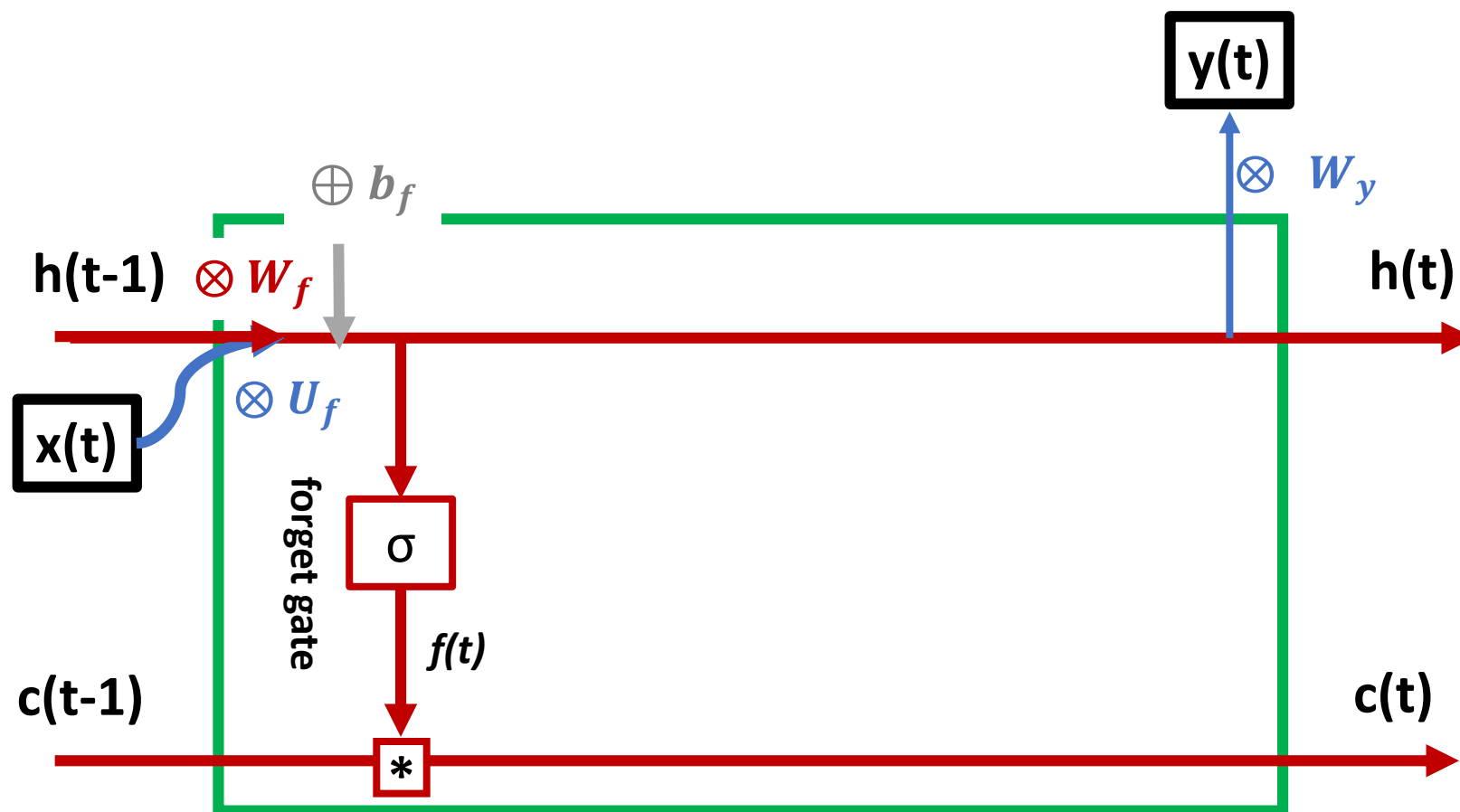
the LSTM cell





$$f(t) = \sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + b_f)$$

* element – wise multiplication

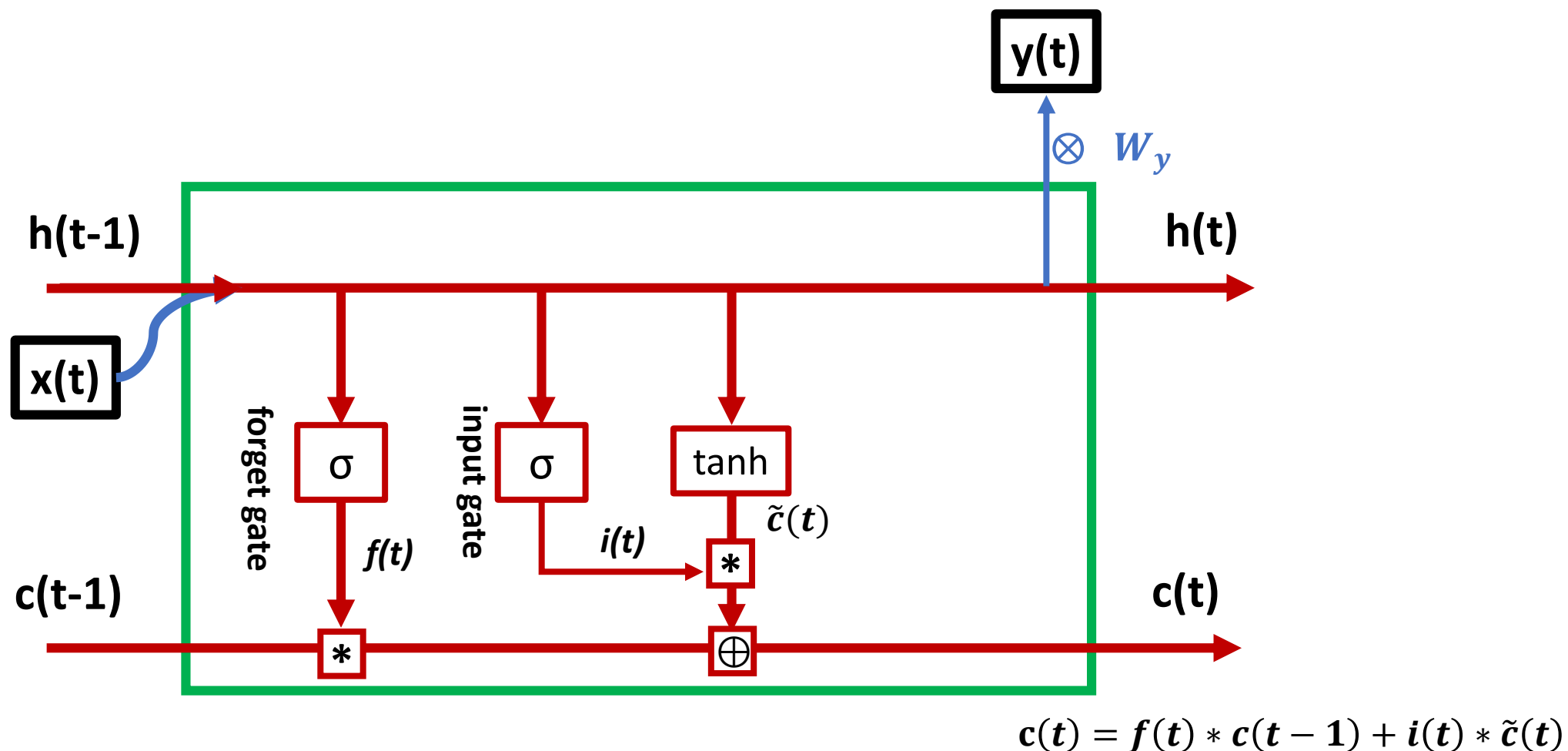




* element – wise multiplication

$$f(t) = \sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + b_f)$$

$$i(t) = \sigma(U_i \oplus x(t) + W_i \oplus h(t-1) + b_i) \quad \tilde{c}(t) = \tanh(U_g \oplus x(t) + W_g \oplus h(t-1) + b_g)$$



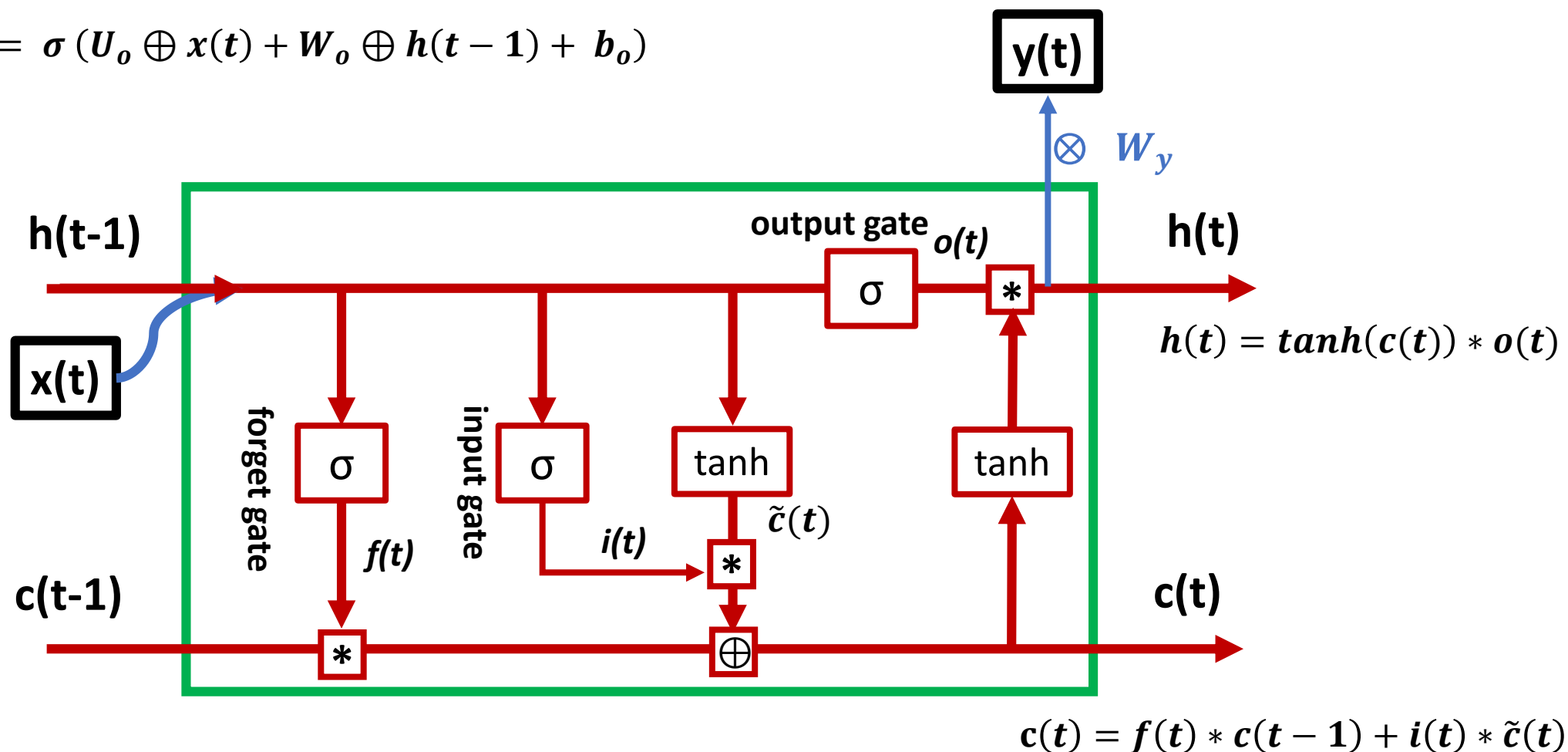


* element – wise multiplication

$$f(t) = \sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + b_f)$$

$$i(t) = \sigma(U_i \oplus x(t) + W_i \oplus h(t-1) + b_i) \quad \tilde{c}(t) = \tanh(U_g \oplus x(t) + W_g \oplus h(t-1) + b_g)$$

$$o(t) = \sigma(U_o \oplus x(t) + W_o \oplus h(t-1) + b_o)$$





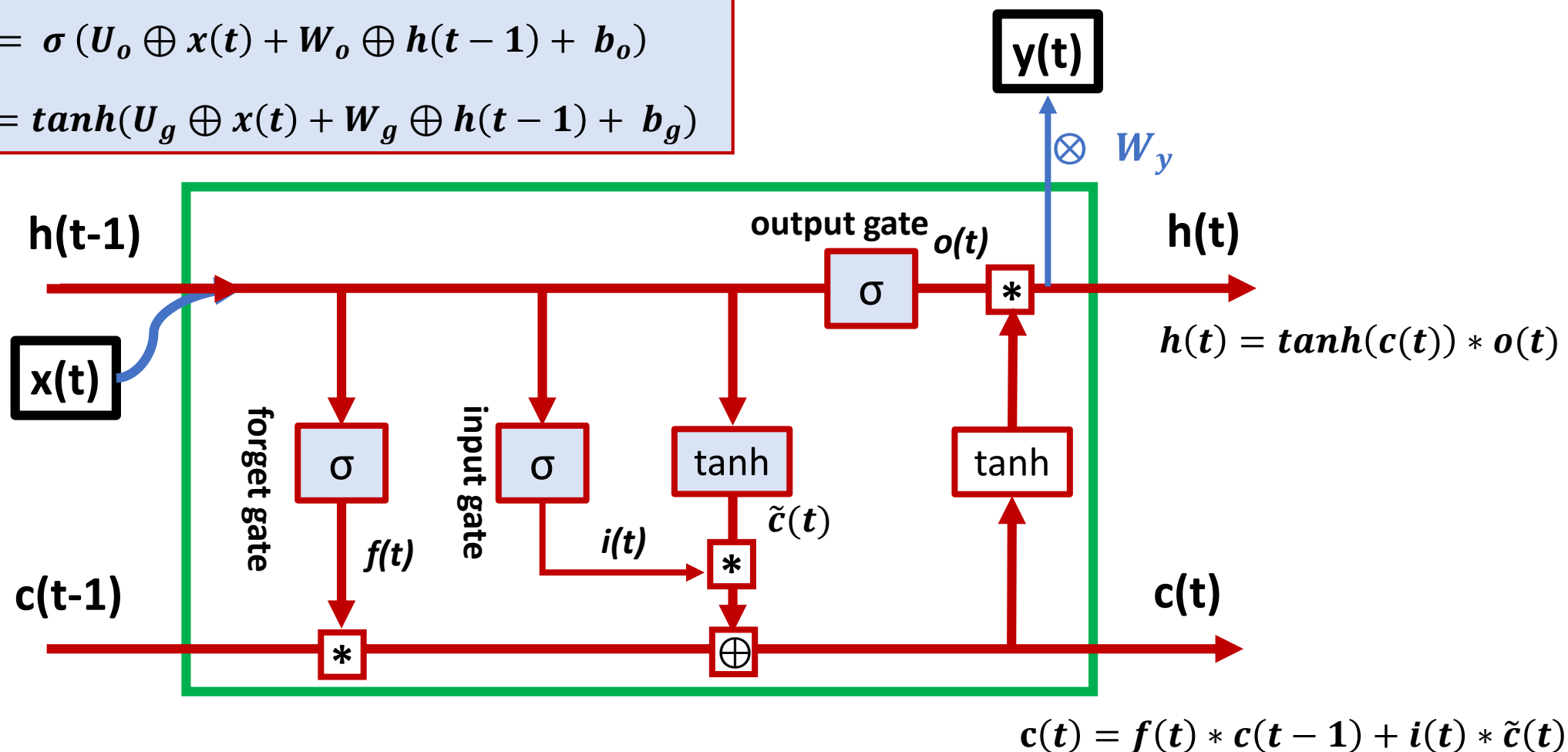
* element – wise multiplication

$$f(t) = \sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + b_f)$$

$$i(t) = \sigma(U_i \oplus x(t) + W_i \oplus h(t-1) + b_i)$$

$$o(t) = \sigma(U_o \oplus x(t) + W_o \oplus h(t-1) + b_o)$$

$$\tilde{c}(t) = \tanh(U_g \oplus x(t) + W_g \oplus h(t-1) + b_g)$$





$$f(t) = \sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + b_f)$$

$$i(t) = \sigma(U_i \oplus x(t) + W_i \oplus h(t-1) + b_i)$$

$$o(t) = \sigma(U_o \oplus x(t) + W_o \oplus h(t-1) + b_o)$$

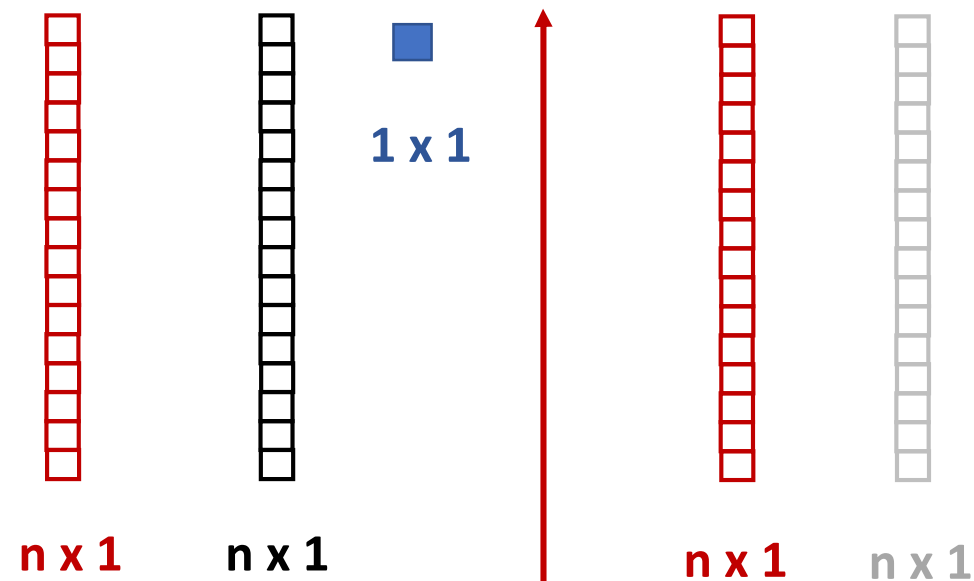
$$\tilde{c}(t) = \tanh(U_g \oplus x(t) + W_g \oplus h(t-1) + b_g)$$

$$h(t) = \tanh(c(t)) * o(t)$$

$$c(t) = f(t) * c(t-1) + i(t) * \tilde{c}(t)$$

* element - wise multiplication

$$f(t) = \sigma(U_f \oplus x(t) + \boxed{W_f} \oplus h(t-1) + b_f)$$



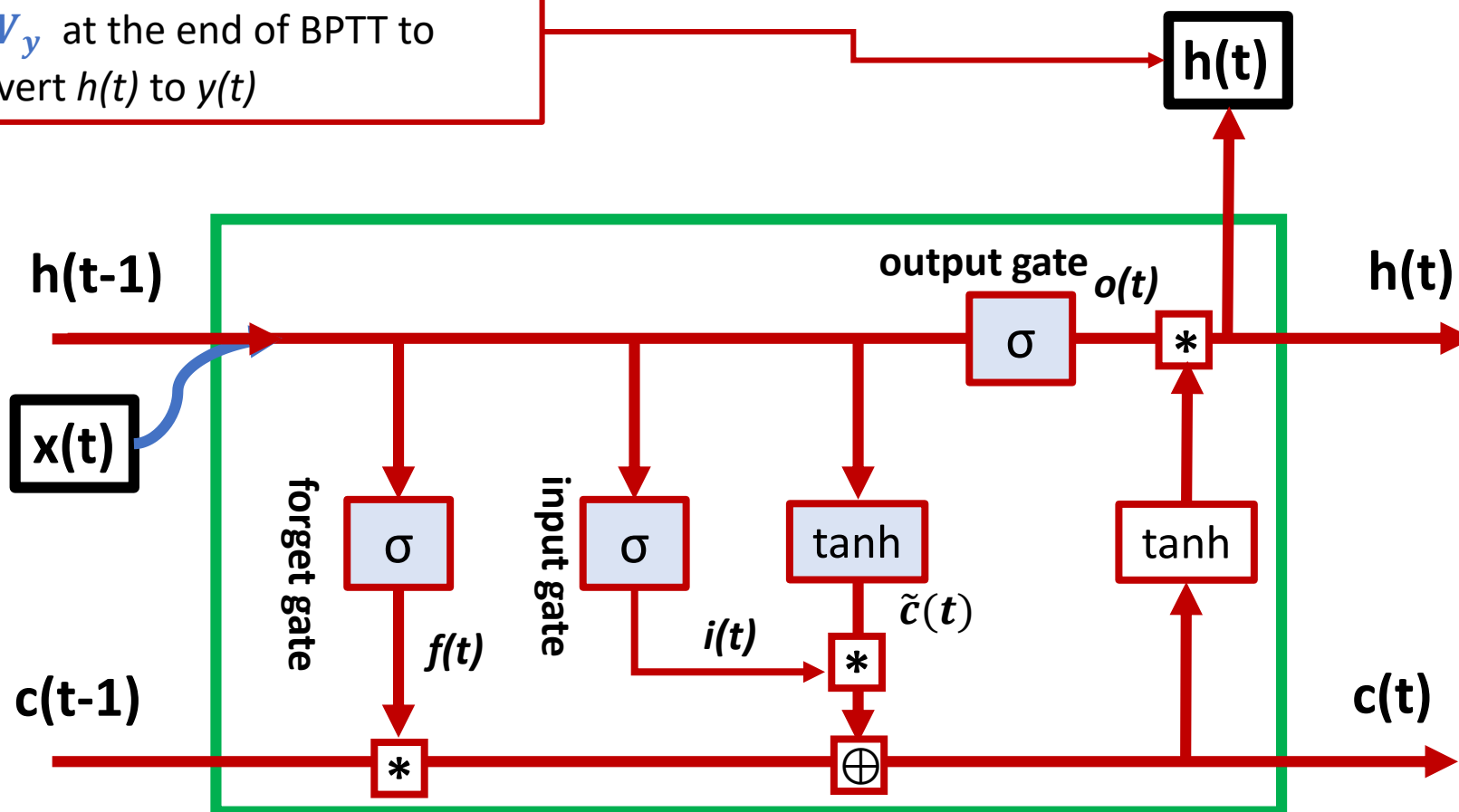
same for $i(t)$, $o(t)$ and $\tilde{c}(t)$



There is one more thing:

* element – wise multiplication

we will add a **dense layer** instead of W_y at the end of BPTT to convert $h(t)$ to $y(t)$





initializing learnables

```
class LSTM():
```

```
    def __init__(self, n_neurons):
```

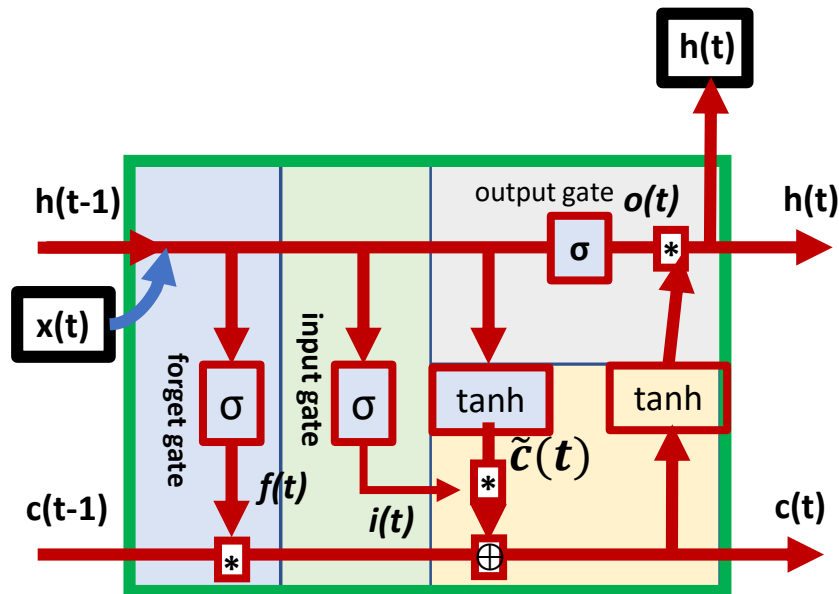
```
        self.n_neurons = n_neurons
```

```
        self.Uf = 0.1*np.random.randn(n_neurons, 1)
        self.bf = 0.1*np.random.randn(n_neurons, 1)
        self.Wf = 0.1*np.random.randn(n_neurons, n_neurons)
```

```
        self.Ui = 0.1*np.random.randn(n_neurons, 1)
        self.bi = 0.1*np.random.randn(n_neurons, 1)
        self.Wi = 0.1*np.random.randn(n_neurons, n_neurons)
```

```
        self.Uo = 0.1*np.random.randn(n_neurons, 1)
        self.bo = 0.1*np.random.randn(n_neurons, 1)
        self.Wo = 0.1*np.random.randn(n_neurons, n_neurons)
```

```
        self.Ug = 0.1*np.random.randn(n_neurons, 1)
        self.bg = 0.1*np.random.randn(n_neurons, 1)
        self.Wg = 0.1*np.random.randn(n_neurons, n_neurons)
```





```
class LSTM():
```

```
    def __init__(self, n_neurons):
```

```
    ...
```

```
    def forward(self, X_t):
```

```
        T = max(X_t.shape)
```

```
        self.T = T
```

```
        self.X_t = X_t
```

```
        n_neurons = self.n_neurons
```

```
        self.H = [np.zeros((n_neurons,1)) for t in range(T+1)]
```

```
        self.C = [np.zeros((n_neurons,1)) for t in range(T+1)]
```

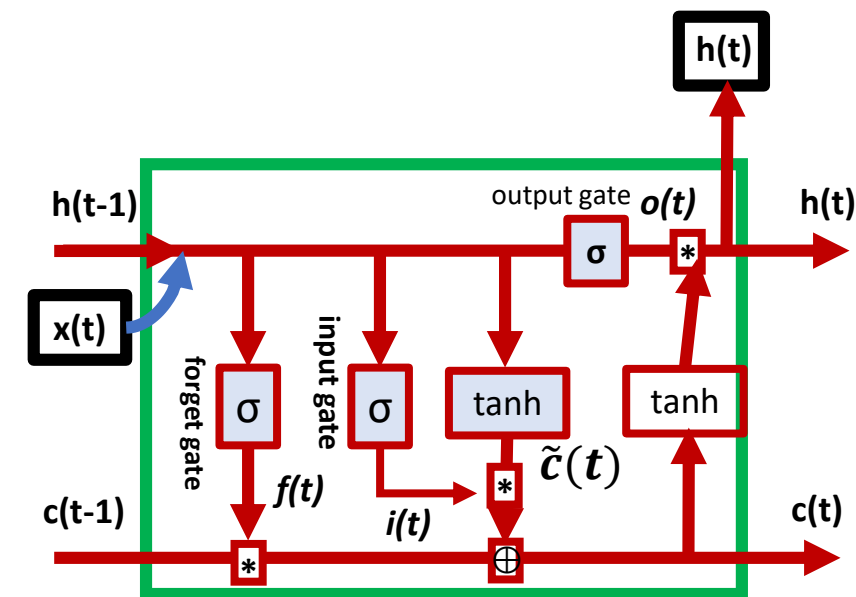
```
        self.C_tilde = [np.zeros((n_neurons,1)) for t in range(T)]
```

```
        self.F = [np.zeros((n_neurons,1)) for t in range(T)]
```

```
        self.O = [np.zeros((n_neurons,1)) for t in range(T)]
```

```
        self.I = [np.zeros((n_neurons,1)) for t in range(T)]
```

keeping track of the input/output of the gates





```
def forward(self, X_t):
```

initializing dweights

```
...
```

```
    self.I = [np.zeros((n_neurons,1)) for t in range(T)]
```

```
    self.dUf = 0.1*np.random.randn(n_neurons, 1)
    self.dbf = 0.1*np.random.randn(n_neurons, 1)
    self.dWf = 0.1*np.random.randn(n_neurons, n_neurons)
```

forget gate

```
    self.dUi = 0.1*np.random.randn(n_neurons, 1)
    self.dbi = 0.1*np.random.randn(n_neurons, 1)
    self.dWi = 0.1*np.random.randn(n_neurons, n_neurons)
```

input gate

```
    self.dUo = 0.1*np.random.randn(n_neurons, 1)
    self.dbo = 0.1*np.random.randn(n_neurons, 1)
    self.dWo = 0.1*np.random.randn(n_neurons, n_neurons)
```

output gate

```
    self.dUg = 0.1*np.random.randn(n_neurons, 1)
    self.dbg = 0.1*np.random.randn(n_neurons, 1)
    self.dWg = 0.1*np.random.randn(n_neurons, n_neurons)
```



```
def forward(self, X_t):
```

instances of activation functions for BPTT

```
...
```

```
    self.dWg = 0.1*np.random.randn(n_neurons, n_neurons)
```

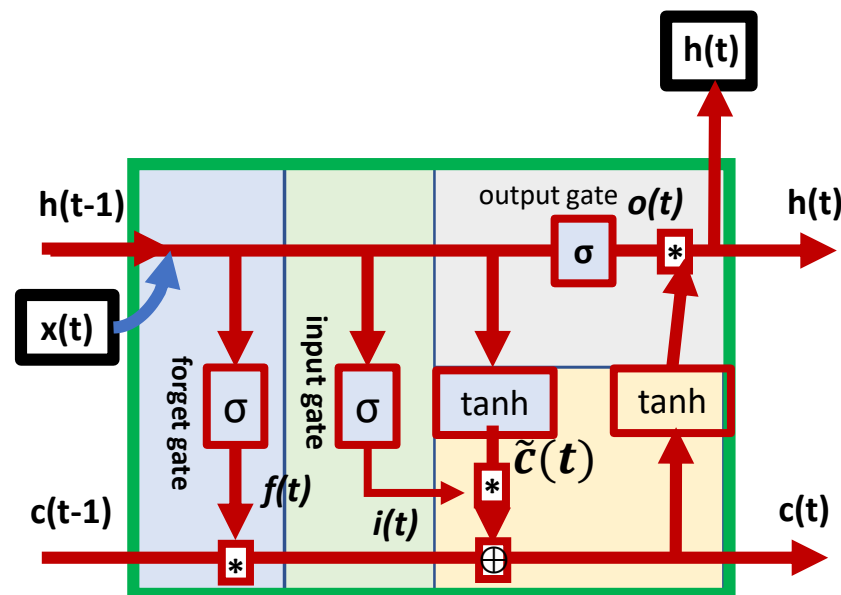
```
    Sigmf = [Sigmoid() for i in range(T)]
```

```
    Sigmi = [Sigmoid() for i in range(T)]
```

```
    Sigmo = [Sigmoid() for i in range(T)]
```

```
    Tanh1 = [Tanh() for i in range(T)]
```

```
    Tanh2 = [Tanh() for i in range(T)]
```





```
def forward(self, X_t):
```

calling the LSTM cell

```
...
```

```
Tanh2 = [Tanh() for i in range(T)]
```

```
ht = self.H[0]
```

```
ct = self.C[0]
```

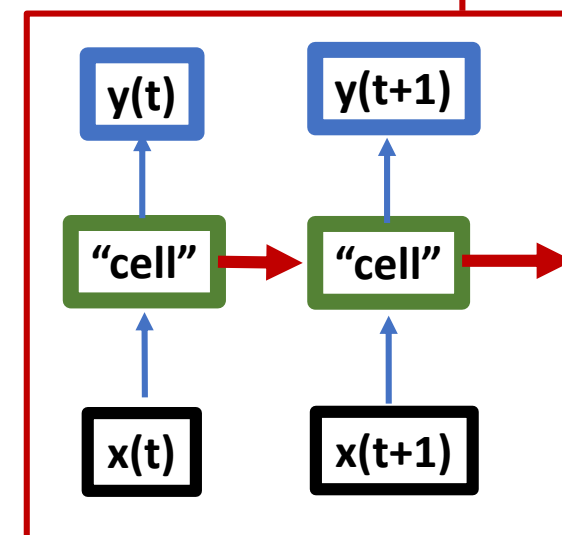
```
[H, C, Sigmf, Sigmi, Sigmo, Tanh1, Tanh2, F, O, I, C_tilde]\n    = self.LSTMCell(X_t, ht, ct, Sigmf, Sigmi, Sigmo, Tanh1, Tanh2,\n                    self.H, self.C, self.F, self.O, self.I, self.C_tilde)
```

```
self.F = F\nself.O = O\nself.I = I\nself.C_tilde = C_tilde
```

```
self.H = H\nself.C = C
```

```
self.Sigmf = Sigmf\nself.Sigmi = Sigmi\nself.Sigmo = Sigmo\nself.Tanh1 = Tanh1\nself.Tanh2 = Tanh2
```

we want to
keep track of
the states
and gates





```
def LSTMCell(self, X_t, ht, ct, Sigmf, Sigmi, Sigmo, Tanh1, Tanh2,\n              H, C, F, O, I, C_tilde):
```

```
    for t, xt in enumerate(X_t):
```

```
        xt = xt.reshape(1,1)
```

forget gate

```
        outf = np.dot(self.Uf, xt) + np.dot(self.Wf, ht) + self.bf
        Sigmf[t].forward(outf)
        ft    = Sigmf[t].output
```

$$f(t) = \sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + b_f)$$

input gate

```
        outi = np.dot(self.Ui, xt) + np.dot(self.Wi, ht) + self.bi
        Sigmi[t].forward(outi)
        it    = Sigmi[t].output
```

$$i(t) = \sigma(U_i \oplus x(t) + W_i \oplus h(t-1) + b_i)$$

output gate

```
        outo = np.dot(self.Uo, xt) + np.dot(self.Wo, ht) + self.bo
        Sigmo[t].forward(outo)
        ot    = Sigmo[t].output
```

$$o(t) = \sigma(U_o \oplus x(t) + W_o \oplus h(t-1) + b_o)$$



```
def LSTMCell(self, X_t, ht, ct, Sigmoid, Sigmoidi, Sigmoido, Tanh1, Tanh2,\n              H, C, F, O, I, C_tilde):\n\n    for t, xt in enumerate(X_t):\n        ...  
        ot = Sigmoid[t].output\n\n        outct_tilde = np.dot(self.Ug, xt) + np.dot(self.Wg, ht) + self.bg  
        Tanh1[t].forward(outct_tilde)  
        ct_tilde = Tanh1[t].output  
         $\tilde{c}(t) = \tanh(U_g \oplus x(t) + W_g \oplus h(t-1) + b_g)$ \n\n        ct = np.multiply(ft, ct) + np.multiply(it, ct_tilde)  
         $c(t) = f(t) * c(t-1) + i(t) * \tilde{c}(t)$ \n\n        Tanh2[t].forward(ct)  
        ht = np.multiply(Tanh2[t].output, ot)  
         $h(t) = \tanh(c(t)) * o(t)$ 
```



```
def LSTMCell(self, X_t, ht, ct, Sigmf, Sigmi, Sigmo, Tanh1, Tanh2,\n              H, C, F, O, I, C_tilde):\n\n    for t, xt in enumerate(X_t):\n        ...\n\n        Tanh2[t].forward(ct)\n        ht = np.multiply(Tanh2[t].output, ot)\n\n        H[t+1]      = ht\n        C[t+1]      = ct\n        C_tilde[t]  = ct_tilde\n\n        F[t]        = ft\n        O[t]        = ot\n        I[t]        = it\n\n    return(H, C, Sigmf, Sigmi, Sigmo, Tanh1, Tanh2, F, O, I, C_tilde)
```

we want to keep track of the states and gates



let us test the code

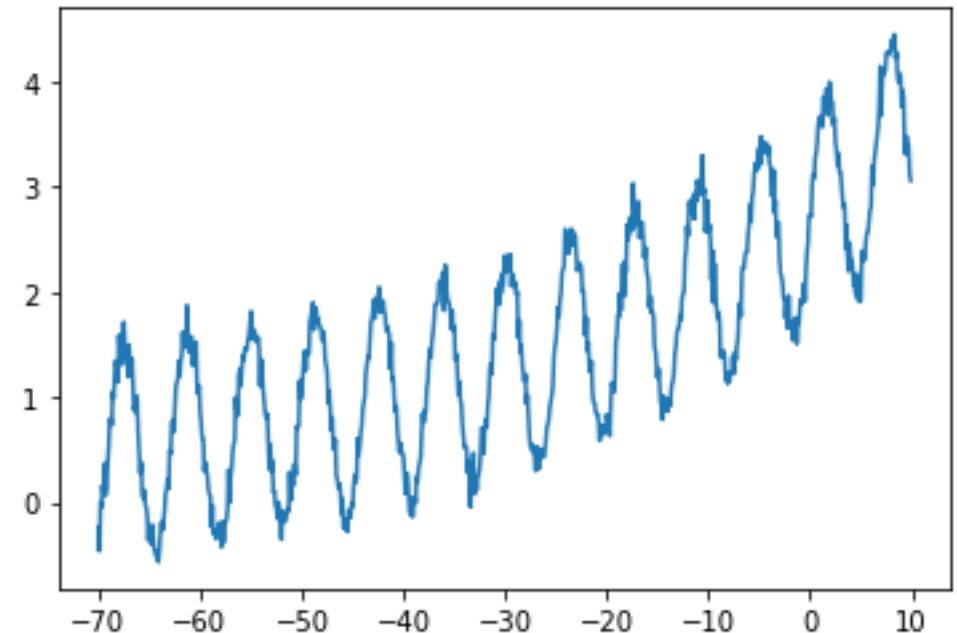
```
import numpy as np
import matplotlib.pyplot as plt

X_t = np.arange(-70,10,0.1)
X_t = X_t.reshape(len(X_t),1)
Y_t = np.sin(X_t) + 0.1*np.random.randn(len(X_t),1) + np.exp((0.5*X_t + 20)*0.05)

plt.plot(X_t, Y_t)
plt.show()
```

aim:

- prediction $Y_t(t) \rightarrow Y_t(t + dt)$



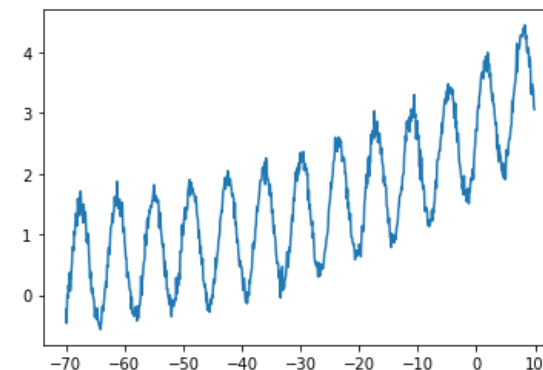


let us test the code

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

X_t = np.arange(-70, 10, 0.1)
X_t = X_t.reshape(len(X_t), 1)
Y_t = np.sin(X_t) + 0.1*np.random.randn(len(X_t), 1) + np.exp((0.5*X_t + 20)*0.05)

plt.plot(X_t, Y_t)
plt.show()
```



```
from LSTM import *
```

```
lstm = LSTM(n_neurons = 200)
lstm.forward(X_t)
```

```
for h in lstm.H:
    plt.plot(np.arange(20), h[0:20], 'k-', linewidth = 1, alpha = 0.05)
```

do the same for F, I and O



let us test the code

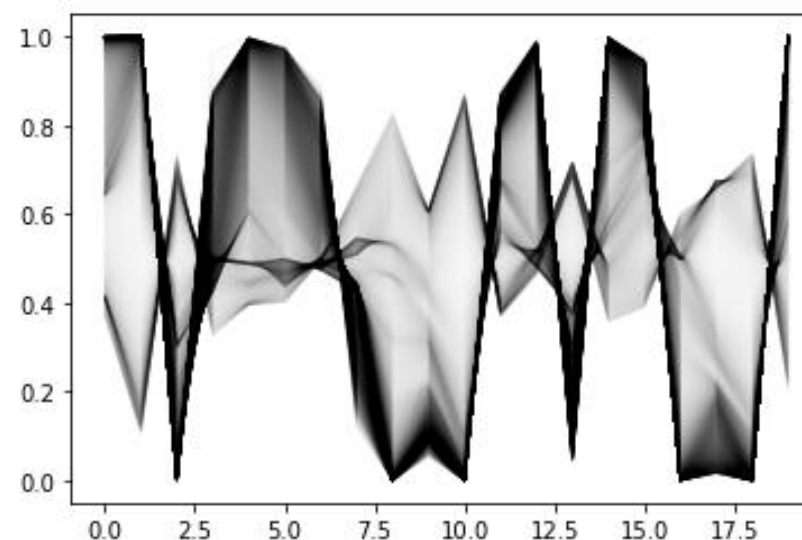
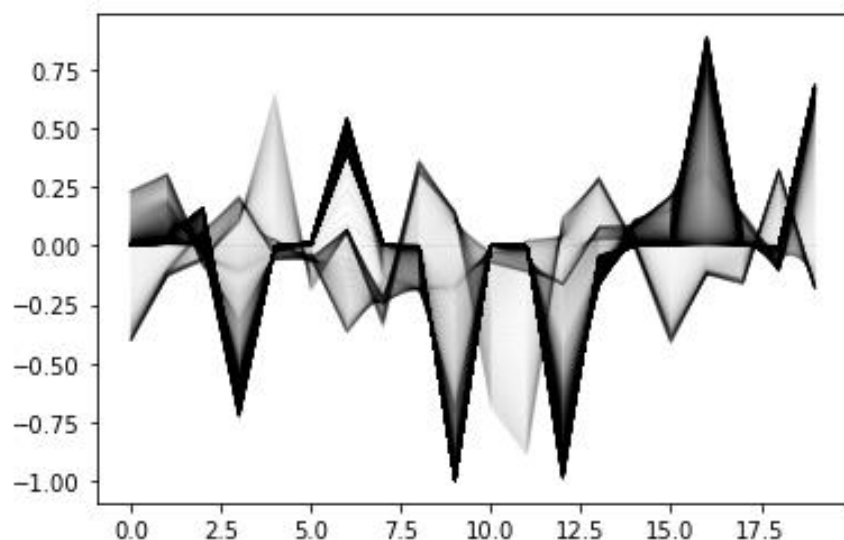
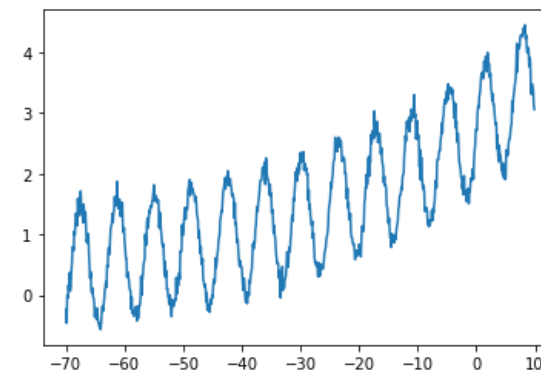
```
from LSTM import *
```

```
lstm = LSTM(n_neurons = 200)
```

```
lstm.forward(X_t)
```

```
for h in lstm.H:
```

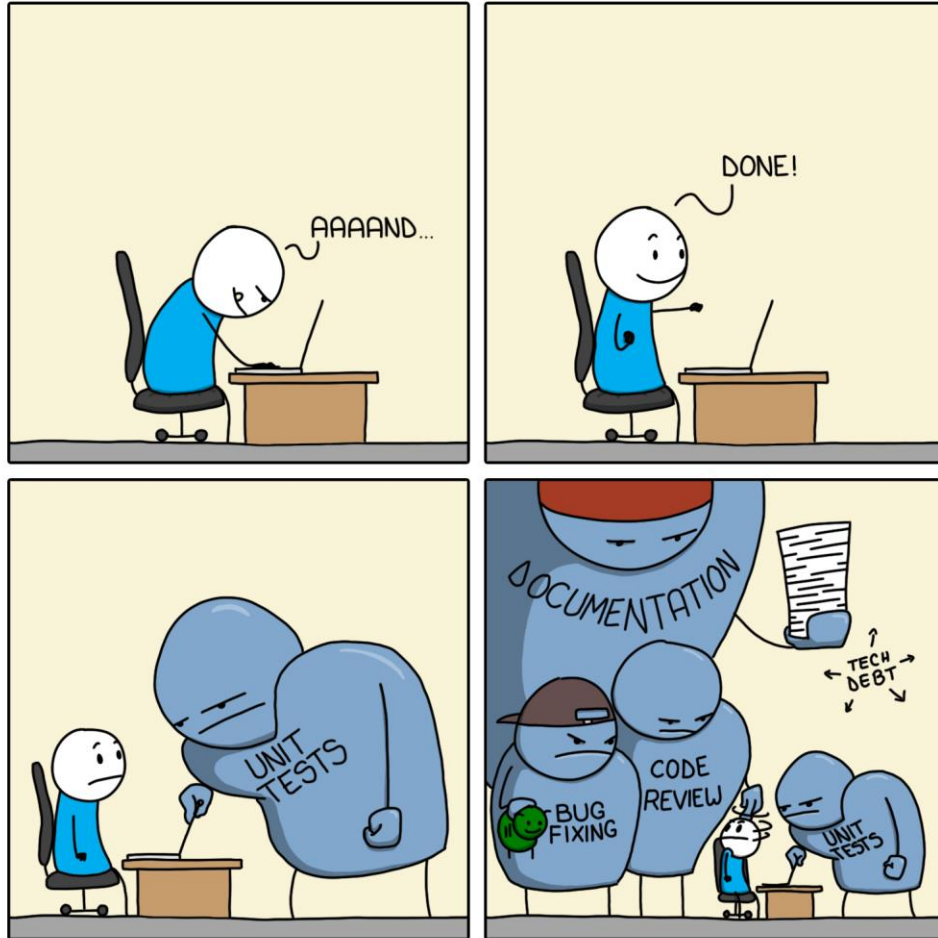
```
    plt.plot(np.arange(20), h[0:20], 'k-', linewidth = 1, alpha = 0.05)
```





FEATURE COMPLETE

MONKEYUSER.COM



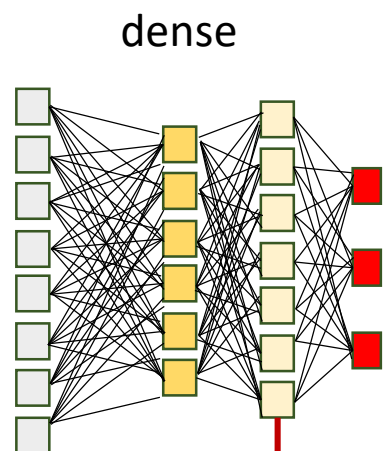
outline:

- *the idea*
- *from a classical RNN cell to a LSTM*
- **BackPropagation Through Time**
- *full backpropagation*
- *modifying the SGD optimizer*
- *running and testing the package*



see lecture
ANN and CNN from Scratch

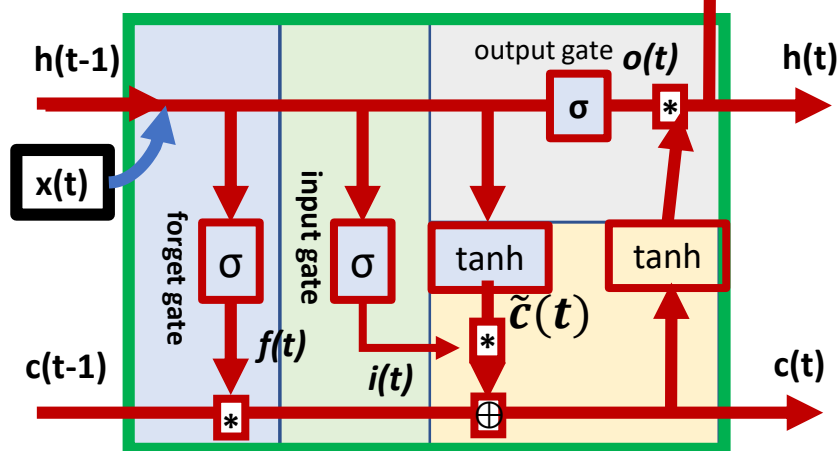
$$L = \frac{1}{2} \sum_{t=1}^T [\hat{y}(t) - y(t)]^2$$



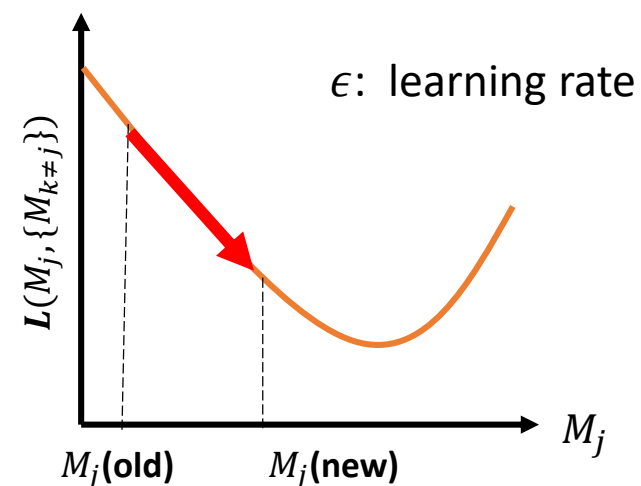
$$dY = dY(U_f, W_f, b_f, U_i, W_i, b_i, U_o, W_o, b_o, U_g, W_g, b_g)$$

dense.dinputs

$h(t)$



$$\Delta M_j = -\epsilon \frac{d L(M_j, \{M_{k \neq j}\})}{d M_j}$$





$$\Delta = \Delta(U_f, W_f, b_f, U_i, W_i, b_i, U_o, W_o, b_o, U_g, W_g, b_g)$$

$$h(t) = \boxed{\tanh(c(t)) * o(t)}$$

$$c(t) = \boxed{f(t) * c(t-1)} + i(t) * \tilde{c}(t)$$

$$f(t) = \boxed{\sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + b_f)}$$

$$\begin{aligned} \frac{dh(t)}{dU_f} &= \frac{\partial h(t)}{\partial \tanh} \frac{\partial \tanh}{\partial c(t)} * o(t) \frac{\partial c(t)}{\partial f(t)} \frac{\partial f(t)}{\partial \sigma} \frac{\partial \sigma}{\partial U_f} \\ &= c(t-1) = x(t) \end{aligned}$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Tanh2[t].backward(dht)
dtanh2 = Tanh2[t].dinputs
```

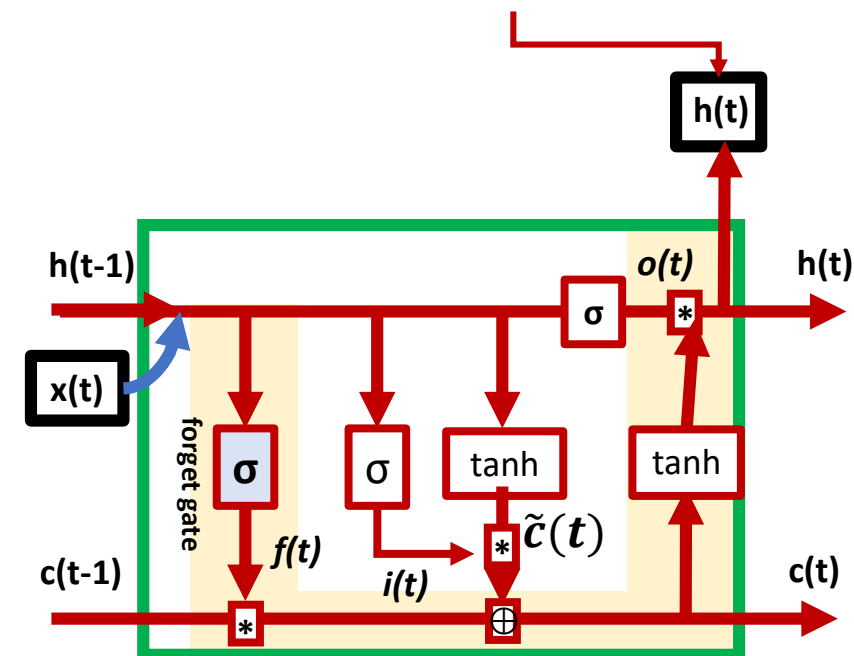
```
dhtdtanh = np.multiply(O[t], dtanh2)
```

```
dctdft = np.multiply(dhtdtanh, C[t-1])
```

```
Sigmf[t].backward(dctdft)
dsigmf = Sigmf[t].dinputs
```

```
dsigmfduf = np.dot(dsigmf, xt)
duf += dsigmfduf
```

$\Delta = \text{dense.dinputs}$





$$\Delta = \Delta (U_f, \mathbf{W}_f, b_f, U_i, W_i, b_i, U_o, W_o, b_o, U_g, W_g, b_g)$$

$$h(t) = \tanh(c(t)) * o(t)$$

$$c(t) = f(t) * c(t-1) + i(t) * \tilde{c}(t)$$

$$f(t) = \sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + b_f)$$

$$\begin{aligned} \frac{dh(t)}{dW_f} &= \frac{\partial h(t)}{\partial \tanh} \frac{\partial \tanh}{\partial c(t)} * o(t) \frac{\partial c(t)}{\partial f(t)} \frac{\partial f(t)}{\partial \sigma} \frac{\partial \sigma}{\partial W_f} \\ &= c(t-1) = h(t-1) \end{aligned}$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Tanh2[t].backward(dht)
dtanh2 = Tanh2[t].dinputs
```

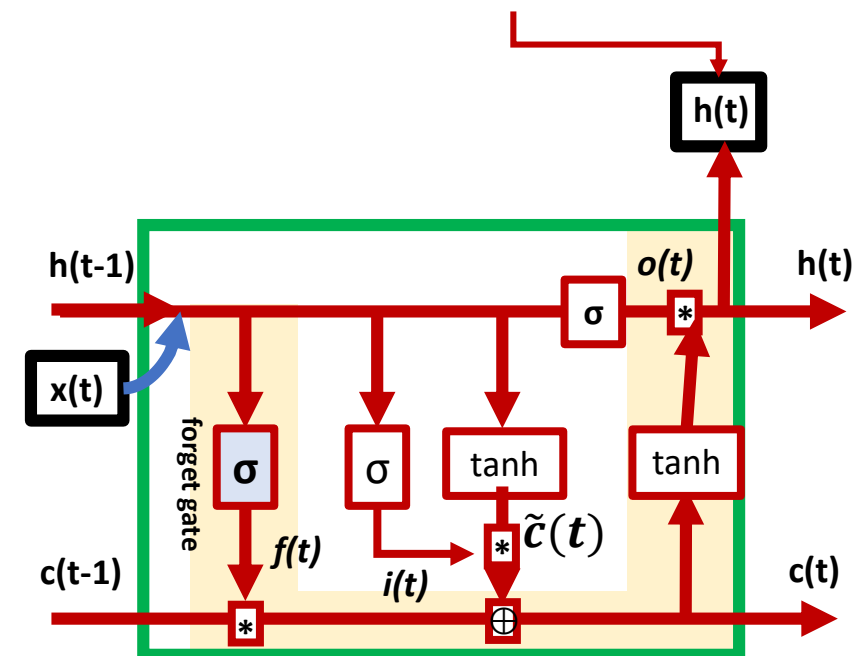
```
dhtdtanh = np.multiply(O[t], dtanh2)
```

```
dctdft = np.multiply(dhtdtanh, C[t-1])
```

```
Sigmf[t].backward(dctdft)
dsigmf = Sigmf[t].dinputs
```

```
dsigmfWf = np.dot(dsigmf, H[t-1].T)
dWf += dsigmfWf
```

$\Delta = \text{dense.dinputs}$





$$\Delta = \Delta (U_f, W_f, \mathbf{b}_f, U_i, W_i, b_i, U_o, W_o, b_o, U_g, W_g, b_g)$$

$$h(t) = \boxed{\tanh(c(t)) * o(t)}$$

$$c(t) = \boxed{f(t) * c(t-1)} + i(t) * \tilde{c}(t)$$

$$\begin{aligned} \frac{dh(t)}{db_f} &= \frac{\partial h(t)}{\partial \tanh} \frac{\partial \tanh}{\partial c(t)} * o(t) \frac{\partial c(t)}{\partial f(t)} \frac{\partial f(t)}{\partial \sigma} \frac{\partial \sigma}{\partial b_f} \\ &= c(t-1) = 1 \end{aligned}$$

$$f(t) = \sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + \boxed{b_f})$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Tanh2[t].backward(dht)
dtanh2 = Tanh2[t].dinputs
```

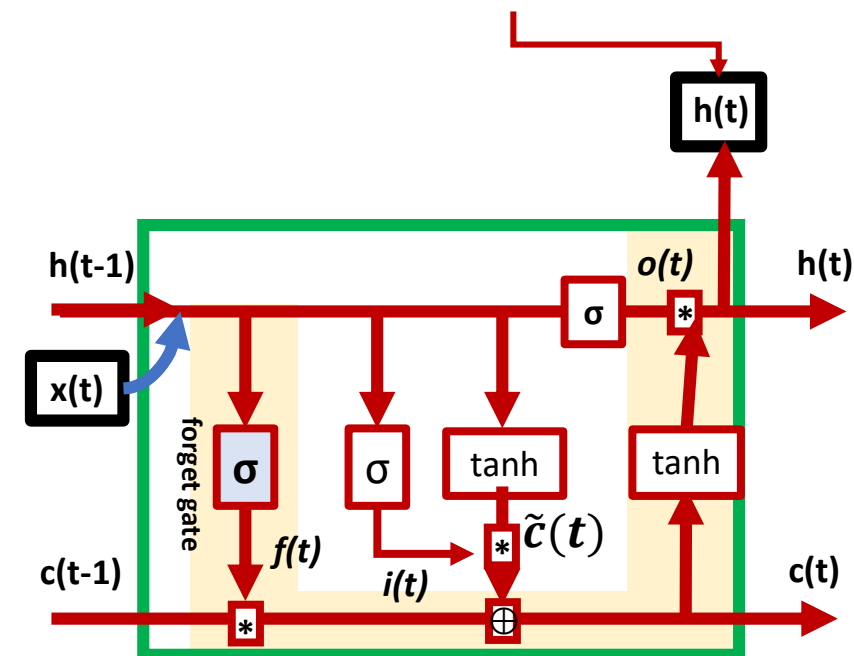
```
dhtdtanh = np.multiply(O[t], dtanh2)
```

```
dctdft = np.multiply(dhtdtanh, C[t-1])
```

```
Sigmf[t].backward(dctdft)
dsigmf = Sigmf[t].dinputs
```

```
dbf += dsigmf
```

$\Delta = \text{dense.dinputs}$





$$\Delta = \Delta(U_f, W_f, b_f, \mathbf{U}_i, W_i, b_i, U_o, W_o, b_o, U_g, W_g, b_g)$$

$$\frac{dh(t)}{dU_i} = \frac{\partial h(t)}{\partial \tanh} \frac{\partial \tanh}{\partial c(t)} * o(t) \frac{\partial c(t)}{\partial i(t)} * \tilde{c}(t) \frac{\partial i(t)}{\partial \sigma} \frac{\partial \sigma}{\partial U_i} = x(t)$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Tanh2[t].backward(dht)
dtanh2 = Tanh2[t].dinputs
```

```
dhtdtanh = np.multiply(O[t], dtanh2)
```

```
dctdit = np.multiply(dhtdtanh, C_tilde[t])
```

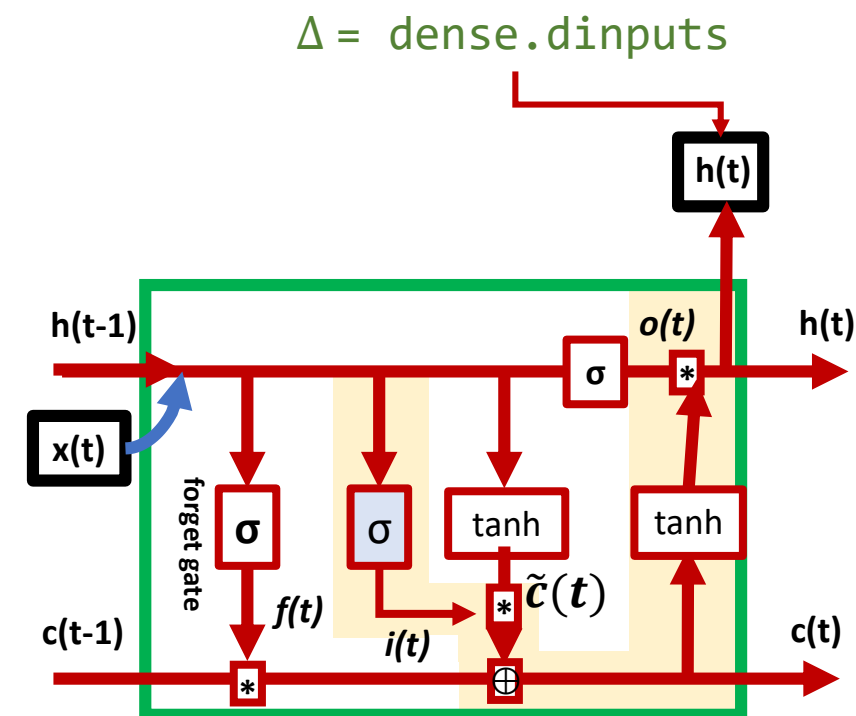
```
Sigmi[t].backward(dctdit)
dsigmi = Sigmi[t].dinputs
```

```
dsigmidUi = np.dot(dsigmi, xt)
dUi += dsigmidUi
```

$$h(t) = \tanh(c(t)) * o(t)$$

$$c(t) = f(t) * c(t-1) + i(t) * \tilde{c}(t)$$

$$i(t) = \sigma(U_i \oplus x(t) + W_i \oplus h(t-1) + b_i)$$





$$\Delta = \Delta (U_f, W_f, b_f, U_i, \mathbf{W}_i, b_i, U_o, W_o, b_o, U_g, W_g, b_g)$$

$$h(t) = \boxed{\tanh(c(t)) * o(t)}$$

$$c(t) = f(t) * c(t-1) + \boxed{i(t) * \tilde{c}(t)}$$

$$i(t) = \sigma(U_i \oplus x(t) + \boxed{W_i \oplus h(t-1)}) + b_i$$

$$\begin{aligned} \frac{dh(t)}{dW_i} &= \frac{\partial h(t)}{\partial \tanh} \frac{\partial \tanh}{\partial c(t)} * o(t) \frac{\partial c(t)}{\partial i(t)} * \tilde{c}(t) \frac{\partial i(t)}{\partial \sigma} \frac{\partial \sigma}{\partial W_i} \\ &= h(t-1) \end{aligned}$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Tanh2[t].backward(dht)
dtanh2 = Tanh2[t].dinputs
```

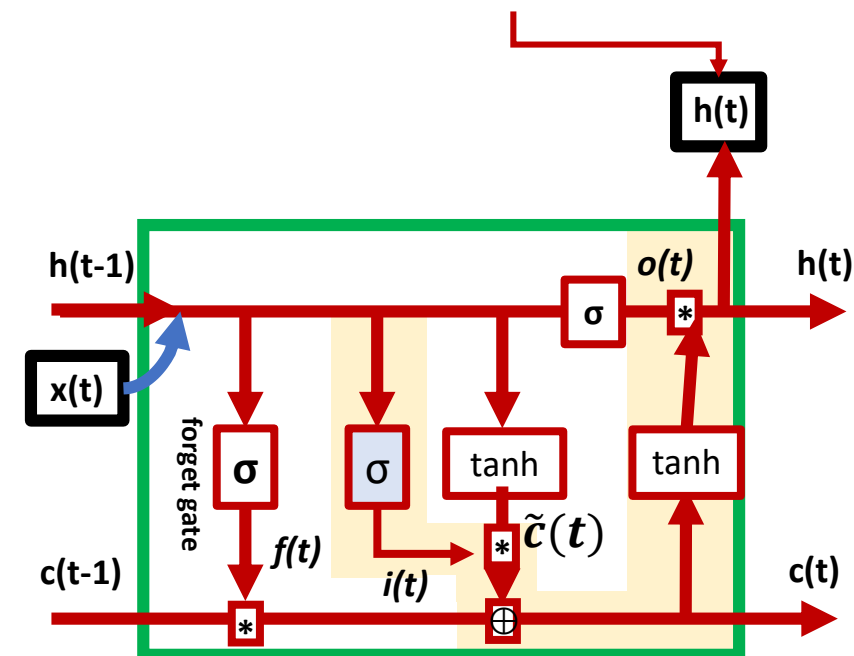
```
dhtdtanh = np.multiply(O[t], dtanh2)
```

```
dctdit = np.multiply(dhtdtanh, C_tilde[t])
```

```
Sigmi[t].backward(dctdit)
dsigmi = Sigmi[t].dinputs
```

```
dsigmidWi = np.dot(dsigmi, H[t-1].T)
dWi += dsigmidWi
```

$\Delta = \text{dense.dinputs}$





$$\Delta = \Delta (U_f, W_f, b_f, U_i, W_i, \mathbf{b}_i, U_o, W_o, b_o, U_g, W_g, b_g)$$

$$\frac{dh(t)}{db_i} = \frac{\partial h(t)}{\partial \tanh} \frac{\partial \tanh}{\partial c(t)} * o(t) \frac{\partial c(t)}{\partial i(t)} * \tilde{c}(t) \frac{\partial i(t)}{\partial \sigma} \frac{\partial \sigma}{\partial b_i} = 1$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Tanh2[t].backward(dht)
dtanh2 = Tanh2[t].dinputs
```

```
dhtdtanh = np.multiply(O[t], dtanh2)
```

```
dctdit = np.multiply(dhtdtanh, C_tilde[t])
```

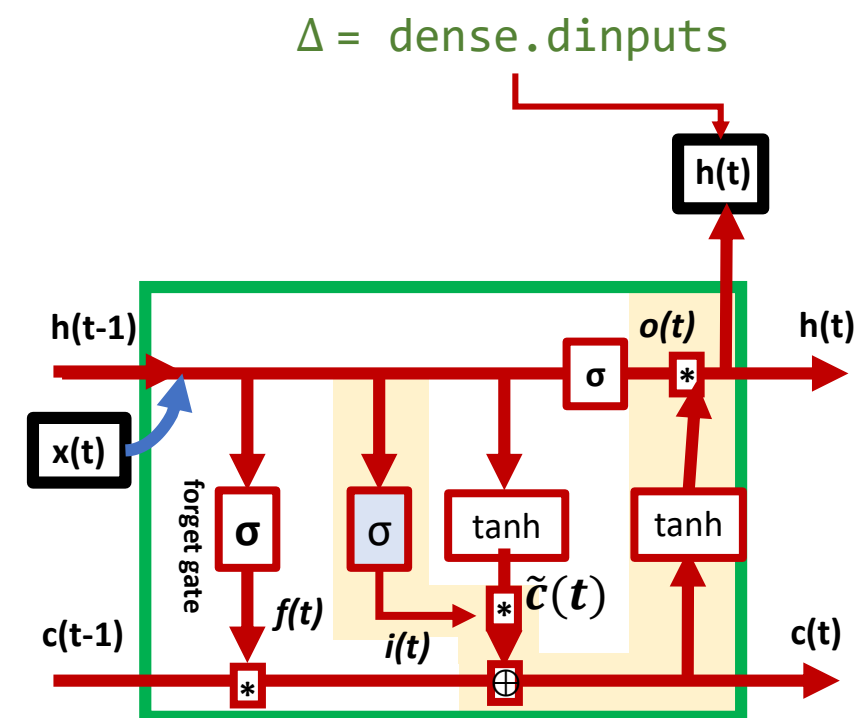
```
Sigmi[t].backward(dctdit)
dsigmi = Sigmi[t].dinputs
```

```
dbi += dsigmi
```

$$h(t) = \tanh(c(t)) * o(t)$$

$$c(t) = f(t) * c(t-1) + i(t) * \tilde{c}(t)$$

$$i(t) = \sigma(U_i \oplus x(t) + W_i \oplus h(t-1) + \mathbf{b}_i)$$





$$\Delta = \Delta (U_f, W_f, b_f, U_i, W_i, b_i, \mathbf{U}_o, W_o, b_o, U_g, W_g, b_g)$$

$$h(t) = \tanh(c(t)) * o(t)$$

$$o(t) = \sigma(U_o \oplus x(t) + W_o \oplus h(t-1) + b_o)$$

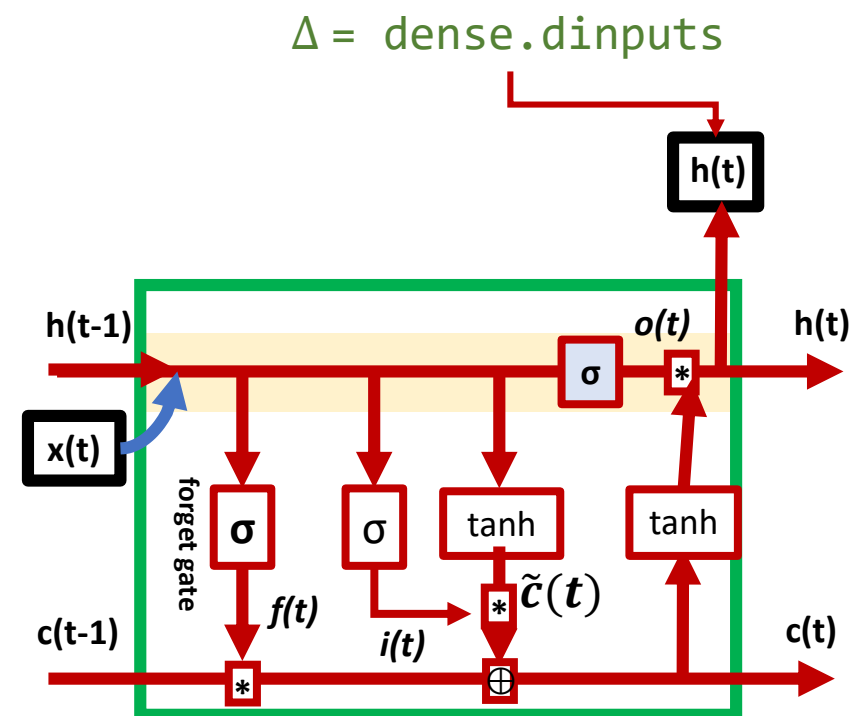
$$\begin{aligned} \frac{dh(t)}{dU_o} &= \frac{\partial h(t)}{\partial o(t)} * \tanh(c(t)) \frac{\partial o(t)}{\partial \sigma} \frac{\partial \sigma}{\partial U_o} \\ &= x(t) \end{aligned}$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Sigmo[t].backward(np.multiply(dht, Tanh2[t].output))
dsigmo = Sigmo[t].dinputs
```

```
dsigmodUo = np.dot(dsigmo, xt)
```

```
dUo += dsigmodUo
```





$$\Delta = \Delta (U_f, W_f, b_f, U_i, W_i, b_i, U_o, \mathbf{W}_o, b_o, U_g, W_g, b_g)$$

$$h(t) = \boxed{\tanh(c(t)) * o(t)}$$

$$o(t) = \sigma (U_o \oplus x(t) + \boxed{W_o \oplus h(t-1)} + b_o)$$

$$\begin{aligned} \frac{dh(t)}{dW_o} &= \frac{\partial h(t)}{\partial o(t)} * \tanh(c(t)) \frac{\partial o(t)}{\partial \sigma} \frac{\partial \sigma}{\partial W_o} \\ &= h(t-1) \end{aligned}$$

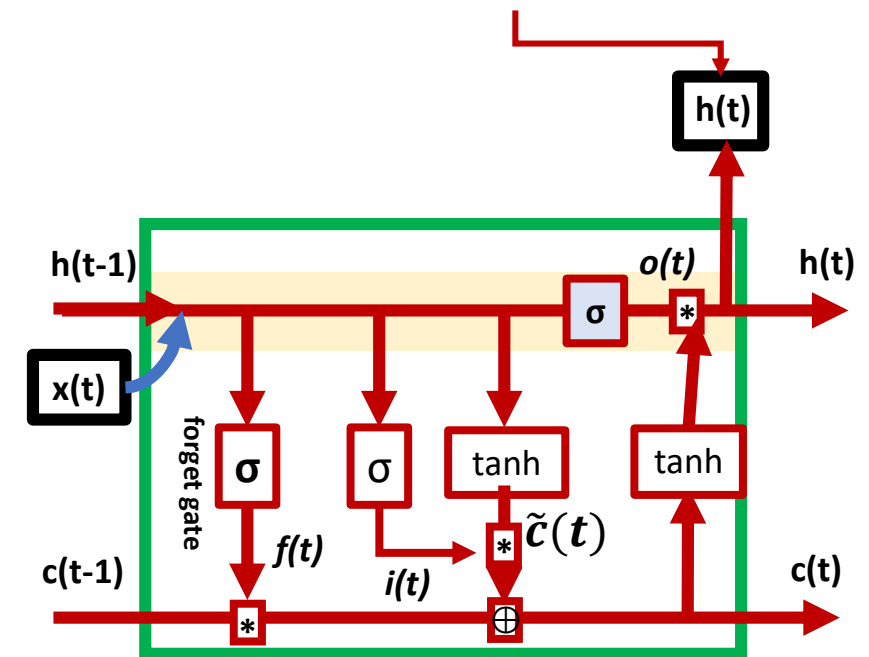
```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Sigmo[t].backward(np.multiply(dht, Tanh2[t].output))
dsigmo = Sigmo[t].dinputs
```

```
dsigmodWo = np.dot(dsigmo, H[t-1].T)
```

```
dWo += dsigmodWo
```

$\Delta = \text{dense.dinputs}$





$$\Delta = \Delta (U_f, W_f, b_f, U_i, W_i, b_i, U_o, W_o, \mathbf{b}_o, U_g, W_g, b_g)$$

$$h(t) = \tanh(c(t)) * o(t)$$

$$o(t) = \sigma (U_o \oplus x(t) + W_o \oplus h(t-1) + \mathbf{b}_o)$$

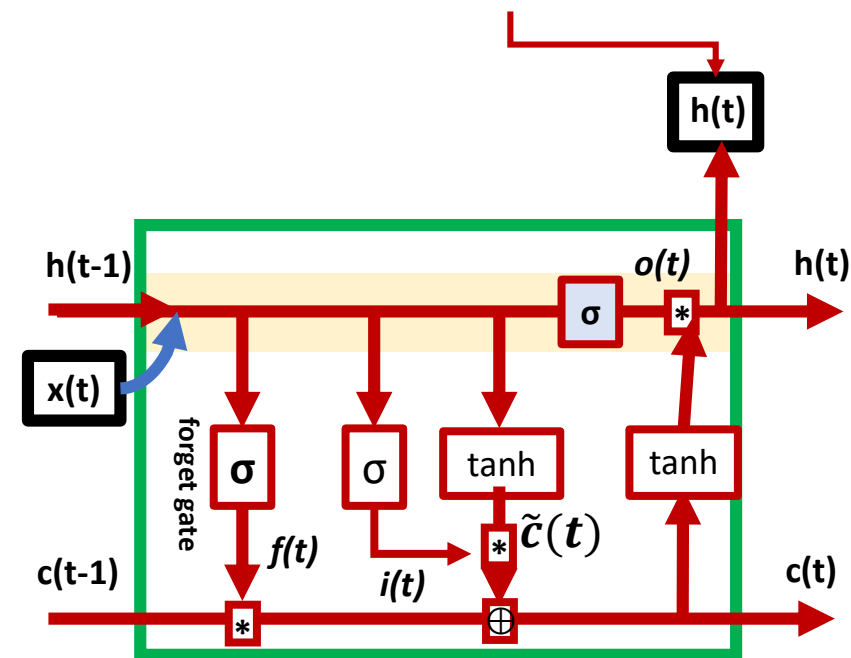
$$\frac{dh(t)}{db_o} = \frac{\partial h(t)}{\partial o(t)} * \tanh(c(t)) \frac{\partial o(t)}{\partial \sigma} \frac{\partial \sigma}{\partial b_o} = 1$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Sigmo[t].backward(np.multiply(dht, Tanh2[t].output))
dsigmo = Sigmo[t].dinputs
```

```
dbo += dsigmo
```

$\Delta = \text{dense.dinputs}$





$$\Delta = \Delta(U_f, W_f, b_f, U_i, W_i, b_i, U_o, W_o, b_o, \mathbf{U}_g, W_g, b_g)$$

$$h(t) = \tanh(c(t)) * o(t)$$

$$c(t) = f(t) * c(t-1) + i(t) * \tilde{c}(t)$$

$$\begin{aligned} \frac{dh(t)}{dU_g} &= \frac{\partial h(t)}{\partial \tanh} \frac{\partial \tanh}{\partial c(t)} * o(t) \frac{\partial c(t)}{\partial \tilde{c}(t)} * i(t) \frac{\partial \tilde{c}(t)}{\partial \tanh} \frac{\partial \tanh}{\partial U_g} \\ &= x(t) \end{aligned}$$

$$\tilde{c}(t) = \tanh(U_g \oplus x(t) + W_g \oplus h(t-1) + b_g)$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Tanh2[t].backward(dht)
dtanh2 = Tanh2[t].dinputs
```

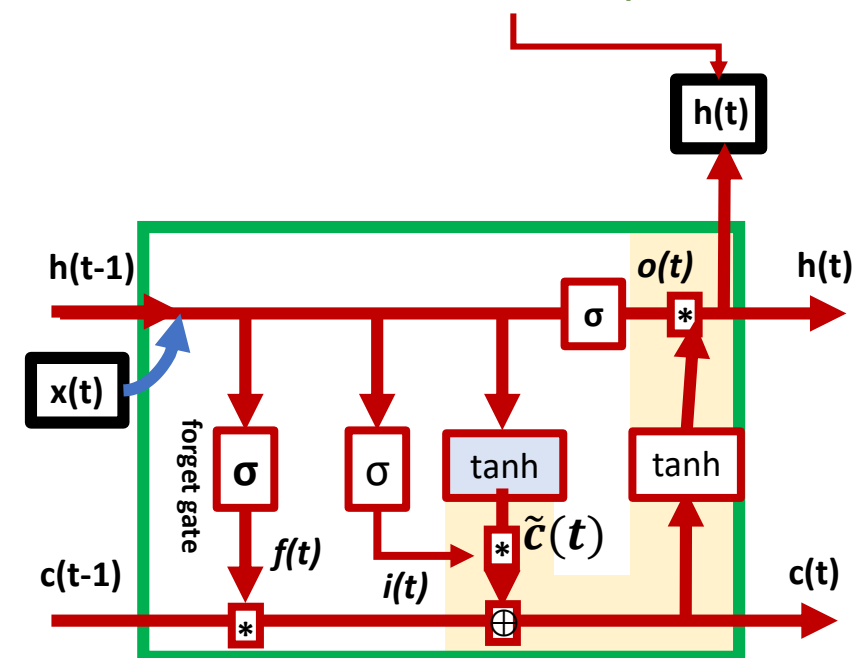
```
dhtdtanh = np.multiply(O[t], dtanh2)
```

```
dctdct_tilde = np.multiply(dhtdtanh, I[t])
```

```
Tanh1[t].backward(dctdct_tilde)
dtanh1 = Tanh1[t].dinputs
```

```
dtanh1dUg = np.dot(dtanh1, xt)
dUg += dtanh1dUg
```

$\Delta = \text{dense.dinputs}$





$$\Delta = \Delta(U_f, W_f, b_f, U_i, W_i, b_i, U_o, W_o, b_o, U_g, \mathbf{W}_g, b_g)$$

$$h(t) = \tanh(c(t)) * o(t)$$

$$c(t) = f(t) * c(t-1) + i(t) * \tilde{c}(t)$$

$$\frac{dh(t)}{dW_g} = \frac{\partial h(t)}{\partial \tanh} \frac{\partial \tanh}{\partial c(t)} * o(t) \frac{\partial c(t)}{\partial \tilde{c}(t)} * i(t) \frac{\partial \tilde{c}(t)}{\partial \tanh} \frac{\partial \tanh}{\partial W_g} \quad \tilde{c}(t) = \tanh(U_g \oplus x(t) + W_g \oplus h(t-1) + b_g)$$

$$= h(t-1)$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Tanh2[t].backward(dht)
dtanh2 = Tanh2[t].dinputs
```

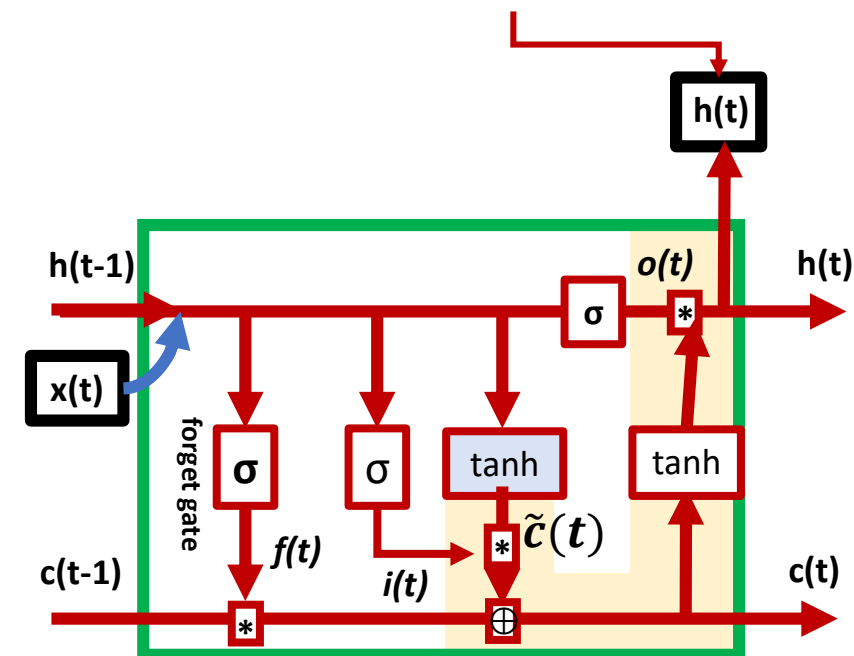
```
dhtdtanh = np.multiply(O[t], dtanh2)
```

```
dctdct_tilde = np.multiply(dhtdtanh, I[t])
```

```
Tanh1[t].backward(dctdct_tilde)
dtanh1 = Tanh1[t].dinputs
```

```
dtanh1dWg = np.dot(dtanh1, H[t-1].T)
dWg += dtanh1dWg
```

$\Delta = \text{dense.dinputs}$





$$\Delta = \Delta(U_f, W_f, b_f, U_i, W_i, b_i, U_o, W_o, b_o, U_g, W_g, \mathbf{b}_g)$$

$$h(t) = \tanh(c(t)) * o(t)$$

$$c(t) = f(t) * c(t-1) + i(t) * \tilde{c}(t)$$

$$\frac{dh(t)}{db_g} = \frac{\partial h(t)}{\partial \tanh} \frac{\partial \tanh}{\partial c(t)} * o(t) \frac{\partial c(t)}{\partial \tilde{c}(t)} * i(t) \frac{\partial \tilde{c}(t)}{\partial \tanh} \frac{\partial \tanh}{\partial b_g} = 1$$

$$\tilde{c}(t) = \tanh(U_g \oplus x(t) + W_g \oplus h(t-1) + b_g)$$

```
dht = dvalues[-1].reshape(self.n_neurons,1)
```

```
Tanh2[t].backward(dht)
dtanh2 = Tanh2[t].dinputs
```

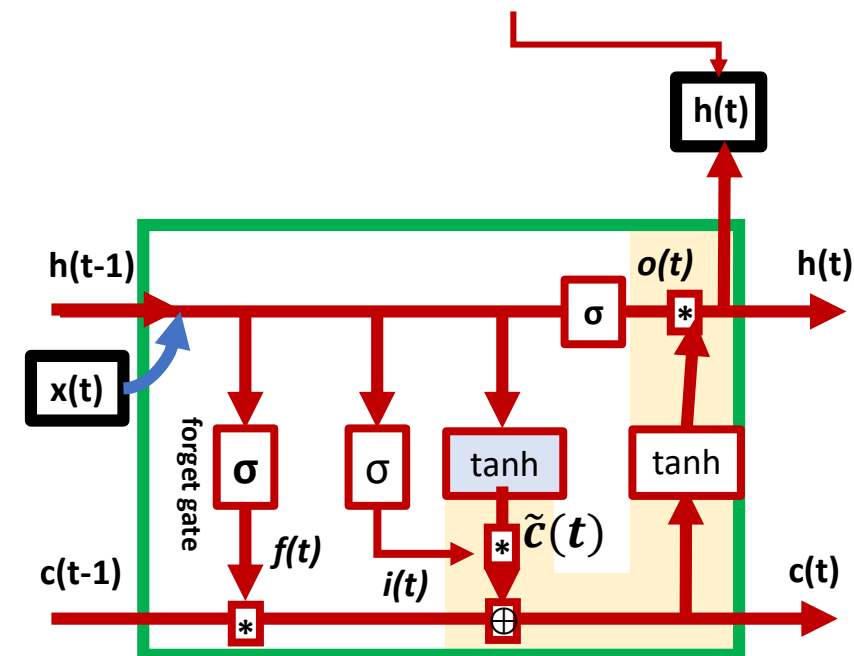
```
dhtdtanh = np.multiply(O[t], dtanh2)
```

```
dctdct_tilde = np.multiply(dhtdtanh, I[t])
```

```
Tanh1[t].backward(dctdct_tilde)
dtanh1 = Tanh1[t].dinputs
```

```
dbg += dtanh1
```

$\Delta = \text{dense.dinputs}$





We finally need $dh(t)$ for the previous cell

$$\frac{dh(t)}{dh(t-1)} =$$

```
dht = np.dot(Wf, dsigmf) + np.dot(Wi, dsigmi) + \
      np.dot(Wo, dsigmo) + np.dot(Wg, dtanh1) + \
      dvalues[t-1,:].reshape(self.n_neurons, 1)
```

$$h(t) = \tanh(c(t)) * o(t)$$

$$c(t) = f(t) * c(t-1) + i(t) * \tilde{c}(t)$$

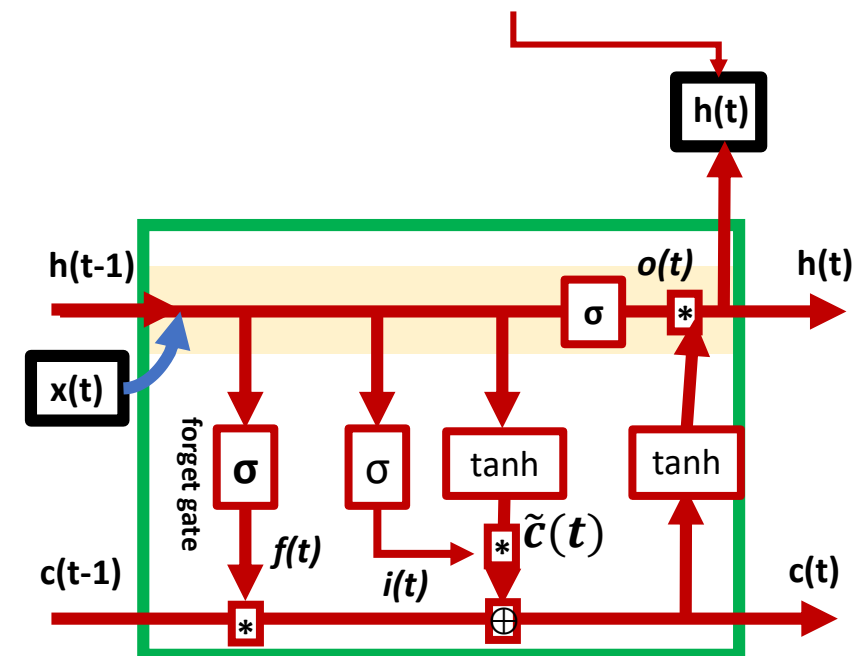
$$f(t) = \sigma(U_f \oplus x(t) + W_f \oplus h(t-1) + b_f)$$

$$\tilde{c}(t) = \tanh(U_g \oplus x(t) + W_g \oplus h(t-1) + b_g)$$

$$i(t) = \sigma(U_i \oplus x(t) + W_i \oplus h(t-1) + b_i)$$

$$o(t) = \sigma(U_o \oplus x(t) + W_o \oplus h(t-1) + b_o)$$

$\Delta = \text{dense.dinputs}$





```
def backward(self, dvalues):  
  
    T      = self.T  
    H      = self.H  
    C      = self.C  
  
    O      = self.O  
    I      = self.I  
    C_tilde = self.C_tilde  
  
    X_t     = self.X_t  
  
    Sigmf   = self.Sigmf  
    Sigmi   = self.Sigmi  
    Sigmoid = self.Sigmoid  
    Tanh1   = self.Tanh1  
    Tanh2   = self.Tanh2
```



```
def backward(self, dvalues):
```

```
    ...
```

```
    Tanh2 = self.Tanh2
```

```
    dht = dvalues[-1,:].reshape(self.n_neurons,1)
```

```
    for t in reversed(range(T)):
```

actual BPTT

```
        xt = X_t[t].reshape(1,1)
```

```
        Tanh2[t].backward(dht)
```

```
        dtanh2 = Tanh2[t].dinputs
```

```
        dhdtanh = np.multiply(0[t], dtanh2)
```

```
        dctdft = np.multiply(dhdtanh, C[t-1])
```

```
        dctdit = np.multiply(dhdtanh, C_tilde[t])
```

```
        dctdct_tilde = np.multiply(dhdtanh, I[t])
```



```
def backward(self, dvalues):
```

```
    ...
```

```
    dht = dvalues[-1].reshape(self.n_neurons, 1)
```

```
    for t in reversed(range(T)):
```

```
        ...
```

```
        dctdct_tilde = np.multiply(dhtdtanh, I[t])
```

```
        Tanh1[t].backward(dctdct_tilde)
```

```
        dtanh1 = Tanh1[t].dinputs
```

```
        Sigmf[t].backward(dctdft)
```

```
        dsigmf = Sigmf[t].dinputs
```

```
        Sigmi[t].backward(dctdit)
```

```
        dsigmi = Sigmi[t].dinputs
```

```
        Sigmot[t].backward(np.multiply(dht, Tanh2[t].output))
```

```
        dsigmo = Sigmot[t].dinputs
```

actual BPTT



```
def backward(self, dvalues):
```

```
    ...  
    dht = dvalues[-1].reshape(self.n_neurons, 1)
```

```
    for t in reversed(range(T)):
```

actual BPTT

```
        ...  
        dsigmo = Sigmoid[t].dinputs  
  
        dsigmfUf = np.dot(dsigmf, xt)  
        dsigmfWf = np.dot(dsigmf, H[t-1].T)  
  
        self.dUf += dsigmfUf  
        self.dWf += dsigmfWf  
        self.dbf += dsigmf  
  
        dsigmiUi = np.dot(dsigmi, xt)  
        dsigmiWi = np.dot(dsigmi, H[t-1].T)  
  
        self.dUi += dsigmiUi  
        self.dWi += dsigmiWi  
        self.dbi += dsigmi
```




```
def backward(self, dvalues):
```

```
    ...  
    dht = dvalues[-1].reshape(self.n_neurons, 1)
```

```
    for t in reversed(range(T)):
```

actual BPTT

```
        ...  
        self.dbi += dsigmi
```

```
        dsigmodUo = np.dot(dsigmo, xt)  
        dsigmodWo = np.dot(dsigmo, H[t-1].T)
```

```
        self.dUo += dsigmodUo  
        self.dWo += dsigmodWo  
        self.dbo += dsigmo
```

```
        dtanh1dUg = np.dot(dtanh1, xt)  
        dtanh1dWg = np.dot(dtanh1, H[t-1].T)
```

```
        self.dUg += dtanh1dUg  
        self.dWg += dtanh1dWg  
        self.dbg += dtanh1
```



```
def backward(self, dvalues):
```

```
    ...  
    dht = dvalues[-1].reshape(self.n_neurons, 1)
```

```
    for t in reversed(range(T)):
```

actual BPTT

```
        ...  
        self.dbg += dtanh1
```

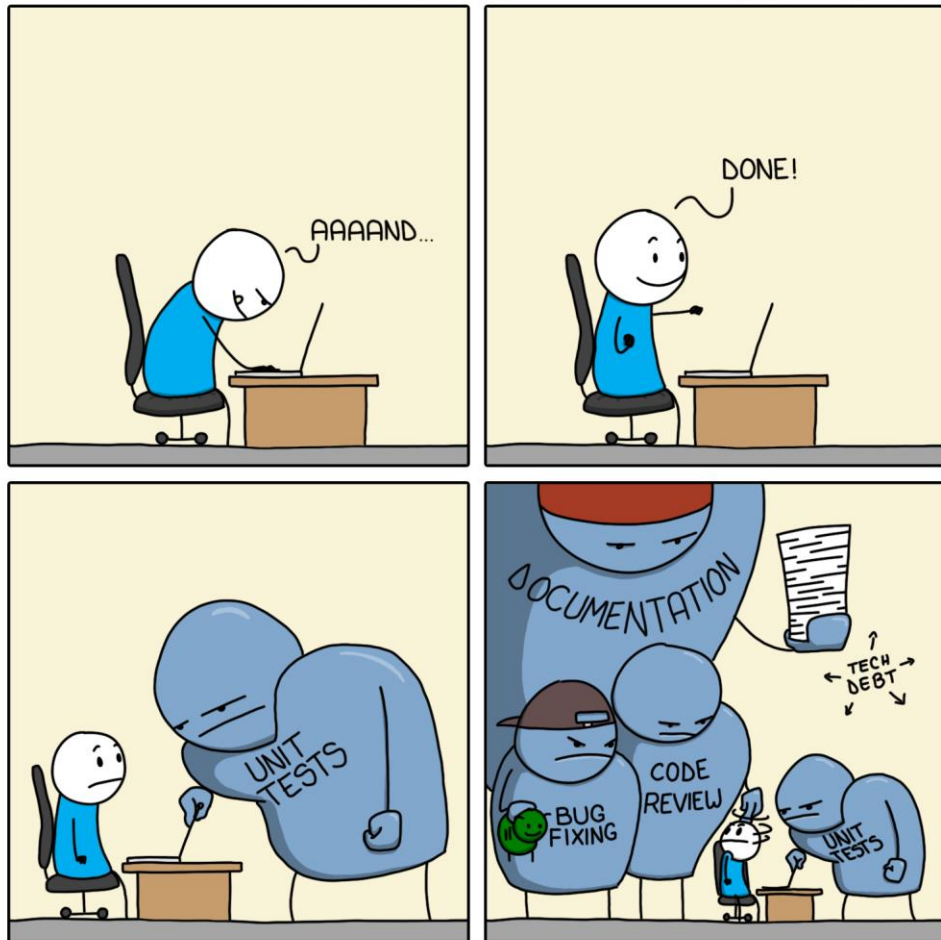
```
        dht = np.dot(self.Wf, dsigmf) + np.dot(self.Wi, dsigmi) +\  
               np.dot(self.Wo, dsigmo) + np.dot(self.Wg, dtanh1) +\  
               dvalues[t-1,:].reshape(self.n_neurons, 1)
```

```
self.H = H
```



FEATURE COMPLETE

MONKEYUSER.COM



outline:

- *the idea*
- *from a classical RNN cell to a LSTM*
- *BackPropagation Through Time*
- **full backpropagation**
- *modifying the SGD optimizer*
- *running and testing the package*



We want to test the backpropagation before we are going to modify our optimizer

```
from LSTM import *
```

```
n_neurons = 200
```

```
lstm = LSTM(n_neurons)
```

```
lstm.forward(X_t)
```

```
T = max(X_t.shape)
```

```
dense1 = Layer_Dense(n_neurons, T)
```

```
dense2 = Layer_Dense(T, 1)
```

```
lr = 1e-5
```

```
Monitor = np.zeros((100))
```

adding a dense layer for

$h(t) \rightarrow \hat{Y}(t)$



We want to test the backpropagation before we are going to modify our optimizer

```
lr = 1e-5
```

```
Monitor = np.zeros((100))
```

```
for i in range(100):
```

```
    lstm.forward(X_t)
```

```
    H = np.array(lstm.H)
```

```
    H = H.reshape((H.shape[0],H.shape[1]))
```

```
    dense1.forward(H[1:,:])
```

```
    dense2.forward(dense1.output)
```

```
    Y_hat = dense2.output
```

```
    dY = Y_hat - Y_t
```

```
    L = float(0.5*np.dot(dY.T,dY)/T)
```

```
    Monitor[i] = L
```



We want to test the backpropagation before we are going to modify our optimizer

```
for i in range(100):  
    ...  
        Monitor[i] = L  
  
        dense2.backward(dY)  
        dense1.backward(dense2.dinputs)  
  
        lstm.backward(dense1.dinputs)  
  
        dense1.weights -= lr*dense1.dweights  
        dense2.weights -= lr*dense2.dweights  
  
        dense1.biases -= lr*dense1.dbiases  
        dense2.biases -= lr*dense2.dbiases
```



We want to test the backpropagation before we are going to modify our optimizer

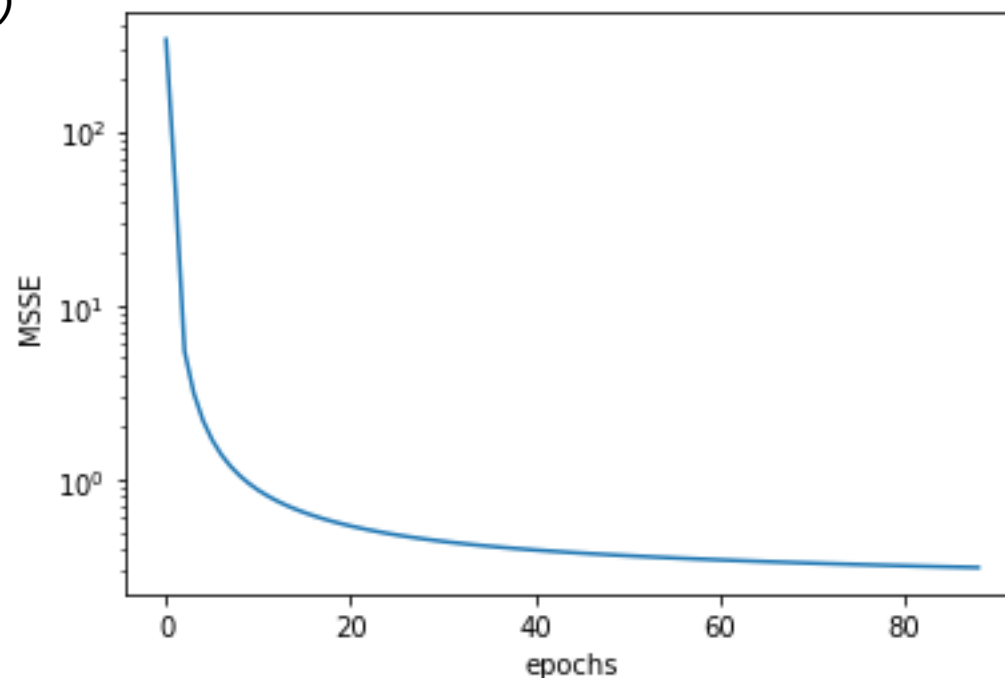
```
for i in range(100):  
    ...  
        dense2.biases -= lr*dense2.dbiases  
  
        lstm.Uf -= lr*lstm.dUf  
        lstm.Ui -= lr*lstm.dUi  
        lstm.Uo -= lr*lstm.dUo  
        lstm.Ug -= lr*lstm.dUg  
  
        lstm.Wf -= lr*lstm.dWf  
        lstm.Wi -= lr*lstm.dWi  
        lstm.Wo -= lr*lstm.dWo  
        lstm.Wg -= lr*lstm.dWg  
  
        lstm.bf -= lr*lstm.dbf  
        lstm.bi -= lr*lstm.dbi  
        lstm.bo -= lr*lstm.dbo  
        lstm.bg -= lr*lstm.dbg  
  
    print(f'current MSSE = {L:.3f}')
```



We want to test the backpropagation before we are going to modify our optimizer

```
for i in range(100):  
    ...  
    print(f'current MSSE = {L:.3f}')
```

```
plt.plot(range(100), Monitor)  
plt.xlabel('epochs')  
plt.ylabel('MSSE')
```

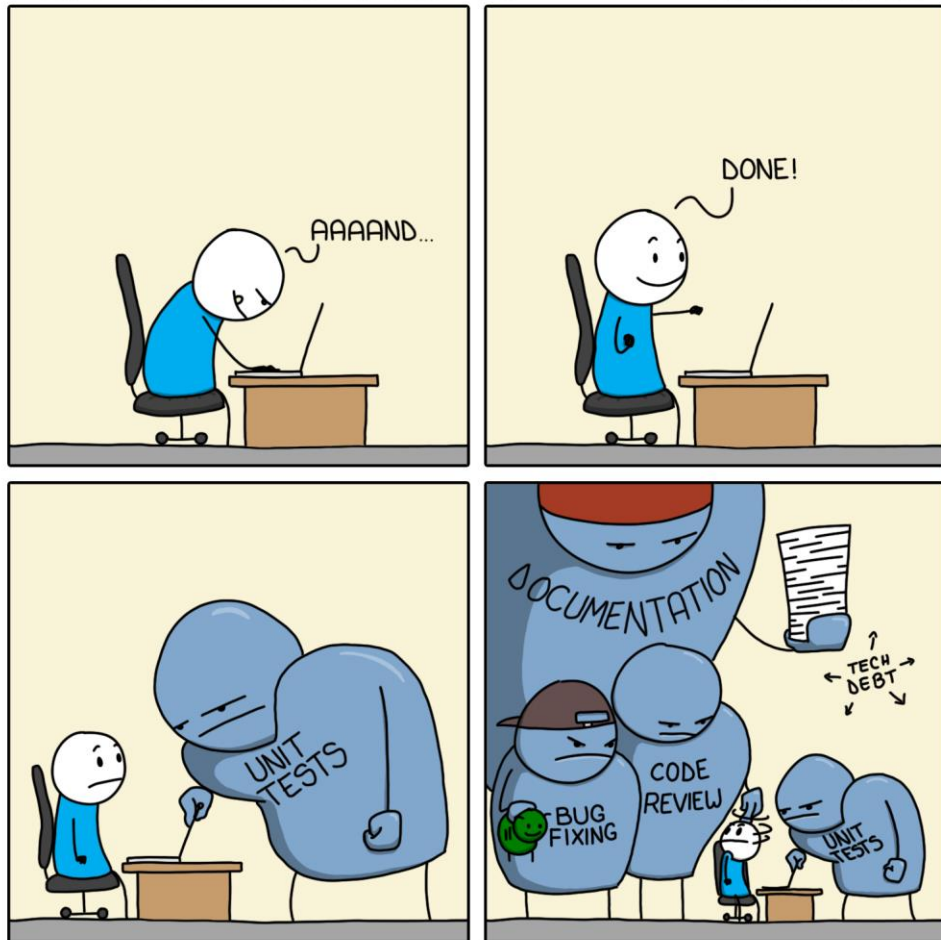


```
current MSSE = 337.707  
current MSSE = 47.685  
current MSSE = 5.465  
current MSSE = 3.140  
current MSSE = 2.186  
current MSSE = 1.688  
current MSSE = 1.388  
current MSSE = 1.191  
current MSSE = 1.052  
current MSSE = 0.949  
current MSSE = 0.870  
current MSSE = 0.807  
current MSSE = 0.756  
current MSSE = 0.713  
current MSSE = 0.677  
current MSSE = 0.647  
current MSSE = 0.620  
current MSSE = 0.597  
current MSSE = 0.576
```




FEATURE COMPLETE

MONKEYUSER.COM



outline:

- *the idea*
- *from a classical RNN cell to a LSTM*
- *BackPropagation Through Time*
- *full backpropagation*
- **modifying the SGD optimizer**
- *running and testing the package*



We want to use **two** optimizer: `Optimizer_SGD` from the previous lecture(s) and...

for more details see
my “ANN from Scratch”
lecture

```
class Optimizer_SGD_LSTM:
```

```
    ...
```

```
    def update_params(self, layer):
```

```
        if self.momentum:
```

```
            if not hasattr(layer, 'Uf_momentums'):
```

```
                layer.Uf_momentums = np.zeros_like(layer.Uf)
```

```
                layer.Ui_momentums = np.zeros_like(layer.Ui)
```

```
                layer.Uo_momentums = np.zeros_like(layer.Uo)
```

```
                layer.Ug_momentums = np.zeros_like(layer.Ug)
```

```
                layer.Wf_momentums = np.zeros_like(layer.Wf)
```

```
                layer.Wi_momentums = np.zeros_like(layer.Wi)
```

```
                layer.Wo_momentums = np.zeros_like(layer.Wo)
```

```
                layer.Wg_momentums = np.zeros_like(layer.Wg)
```

```
                layer.bf_momentums = np.zeros_like(layer.bf)
```

```
                layer.bi_momentums = np.zeros_like(layer.bi)
```

```
                layer.bo_momentums = np.zeros_like(layer.bo)
```

```
                layer.bg_momentums = np.zeros_like(layer.bg)
```

note: try

```
dir(lstm)
```

as alternative



We want to use **two** optimizer: `Optimizer_SGD` from the previous lecture(s) and...

```
class Optimizer_SGD_LSTM:
```

```
...
```

```
def update_params(self, layer):
```

```
    if self.momentum:
```

```
        if not hasattr(layer, 'Uf_momentums'):
```

```
            ...
```

```
            Uf_updates = self.momentum * layer.Uf_momentums - \
                          self.current_learning_rate * layer.dUf
            layer.Uf_momentums = Uf_updates
```

for more details see
my “ANN from Scratch”
lecture

do this for all
learnables of the LSTM



We want to use **two** optimizer: `Optimizer_SGD` from the previous lecture(s) and...

```
class Optimizer_SGD_LSTM:
```

```
...
```

```
def update_params(self, layer):
```

```
    if self.momentum:
```

```
        if not hasattr(layer, 'Uf_momentums'):
```

```
            ...
```

```
            Uf_updates = self.momentum * layer.Uf_momentums - \
                        self.current_learning_rate * layer.dUf
            layer.Uf_momentums = Uf_updates
```

```
            ...
```

```
        else:
```

```
            Uf_updates = -self.current_learning_rate * layer.dUf
            Ui_updates = -self.current_learning_rate * layer.dUi
            Uo_updates = -self.current_learning_rate * layer.dUo
            Ug_updates = -self.current_learning_rate * layer.dUg
```

for more details see
my “ANN from Scratch”
lecture

do this for all
learnables of the LSTM



We want to use **two** optimizer: `Optimizer_SGD` from the previous lecture(s) and...

```
class Optimizer_SGD_LSTM:
```

```
...
```

```
def update_params(self, layer):
```

```
...
```

```
    else:
```

```
        Uf_updates = -self.current_learning_rate * layer.dUf
        Ui_updates = -self.current_learning_rate * layer.dUi
        Uo_updates = -self.current_learning_rate * layer.dUo
        Ug_updates = -self.current_learning_rate * layer.dUg
```

```
        ...
```

```
        layer.Uf += Uf_updates
        layer.Ui += Ui_updates
        layer.Uo += Uo_updates
        layer.Ug += Ug_updates
```

```
        ...
```

```
        lstm.Uf -= lr*lstm.dUf
        lstm.Ui -= lr*lstm.dUi
        lstm.Uo -= lr*lstm.dUo
        lstm.Ug -= lr*lstm.dUg
```

for more details see
my “ANN from Scratch”
lecture

do this for all
learnables of the LSTM



```
optimizerLSTM = Optimizer_SGD_LSTM(learning_rate = 1e-5)
optimizer      = Optimizer_SGD(learning_rate = 1e-5)
Monitor = np.zeros((100))
```

```
for i in range(100):
```

```
    lstm.forward(X_t)
```

```
    H = np.array(lstm.H)
```

```
    H = H.reshape((H.shape[0],H.shape[1]))
```

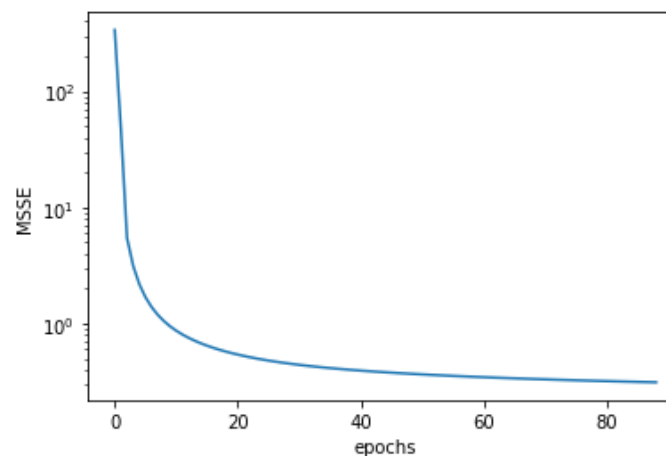
```
    dense1.forward(H[1:,:])
```

```
    dense2.forward(dense1.output)
```

```
    ...
```



```
...  
dense2.forward(dense1.output)  
...  
optimizer_lstm.pre_update_params()  
optimizer.pre_update_params()  
  
optimizer.update_params(dense1)  
optimizer.update_params(dense2)  
  
optimizer_lstm.update_params(lstm)  
  
optimizer_lstm.post_update_params()  
optimizer.post_update_params()
```

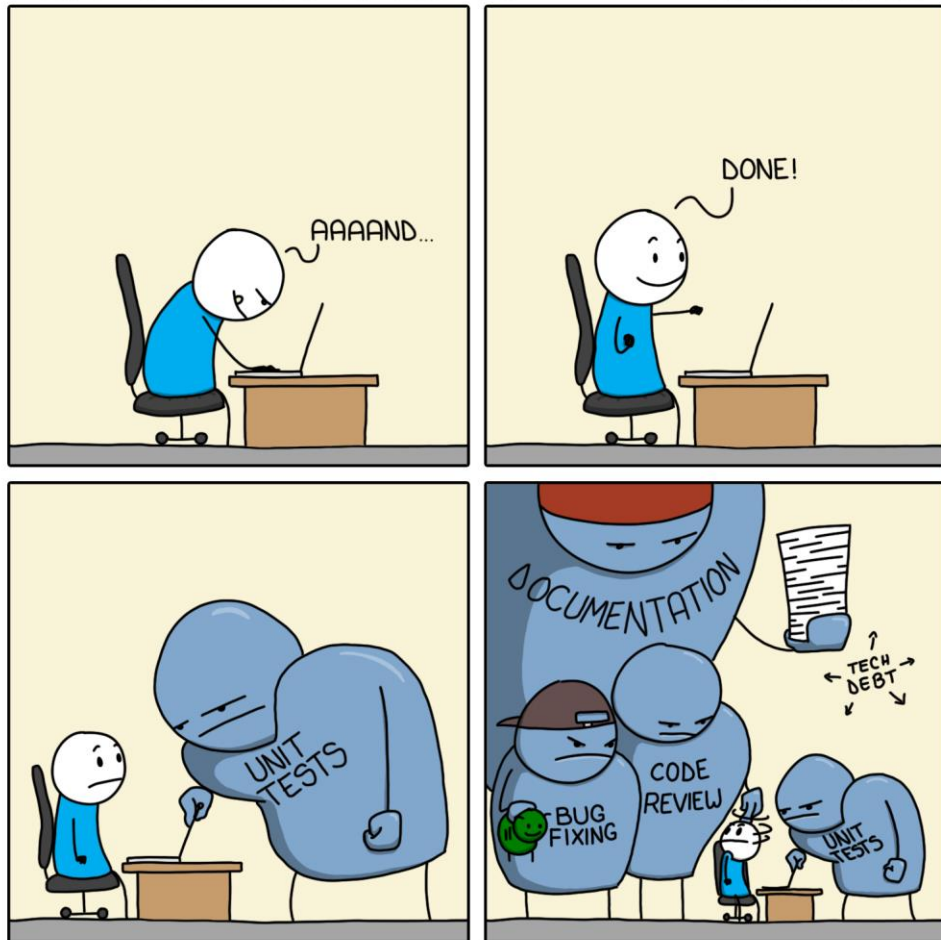


```
current MSSE = 337.707  
current MSSE = 47.685  
current MSSE = 5.465  
current MSSE = 3.140  
current MSSE = 2.186  
current MSSE = 1.688  
current MSSE = 1.388  
current MSSE = 1.191  
current MSSE = 1.052  
current MSSE = 0.949  
current MSSE = 0.870  
current MSSE = 0.807  
current MSSE = 0.756  
current MSSE = 0.713  
current MSSE = 0.677  
current MSSE = 0.647  
current MSSE = 0.620  
current MSSE = 0.597  
current MSSE = 0.576
```



FEATURE COMPLETE

MONKEYUSER.COM

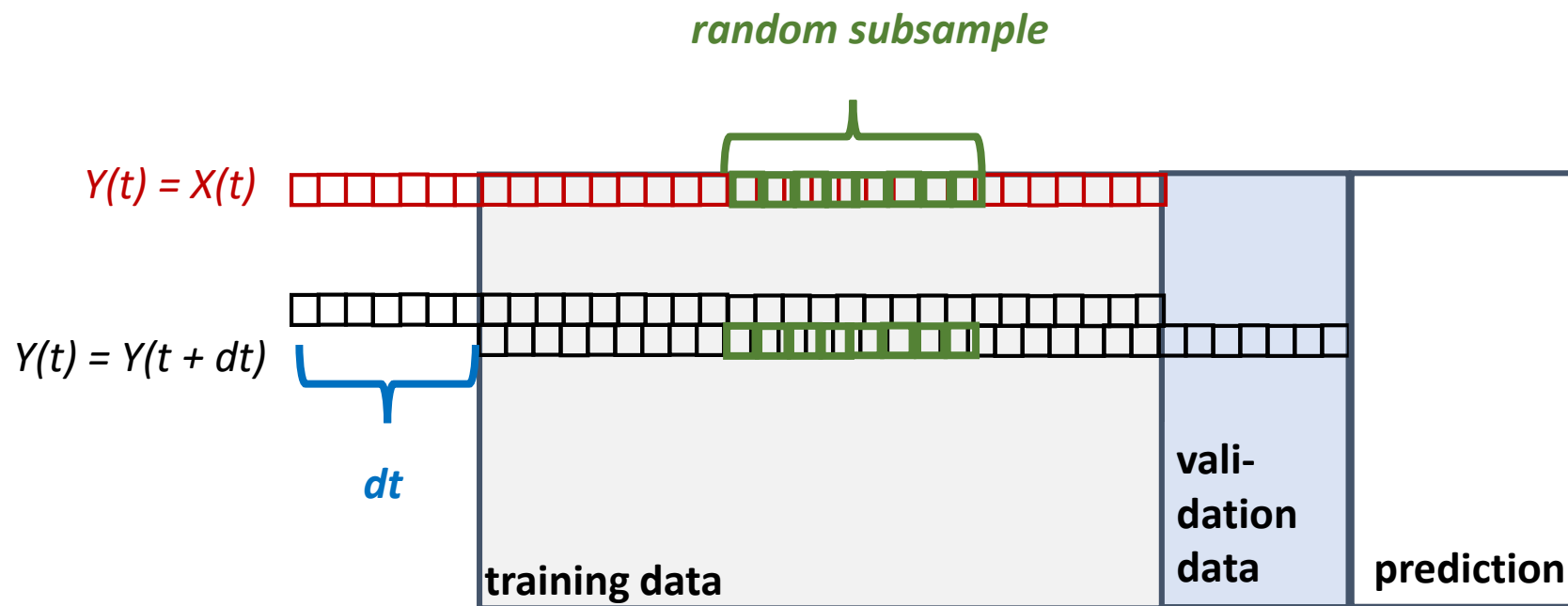


outline:

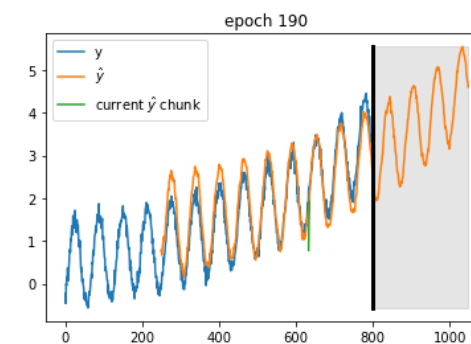
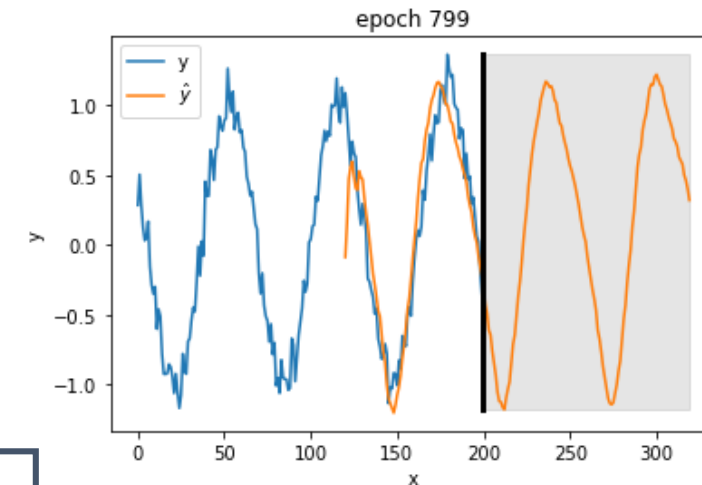
- *the idea*
- *from a classical RNN cell to a LSTM*
- *BackPropagation Through Time*
- *full backpropagation*
- *modifying the SGD optimizer*
- **running and testing the package**



now: $Y(t) \rightarrow Y(t + dt)$



classical RNN

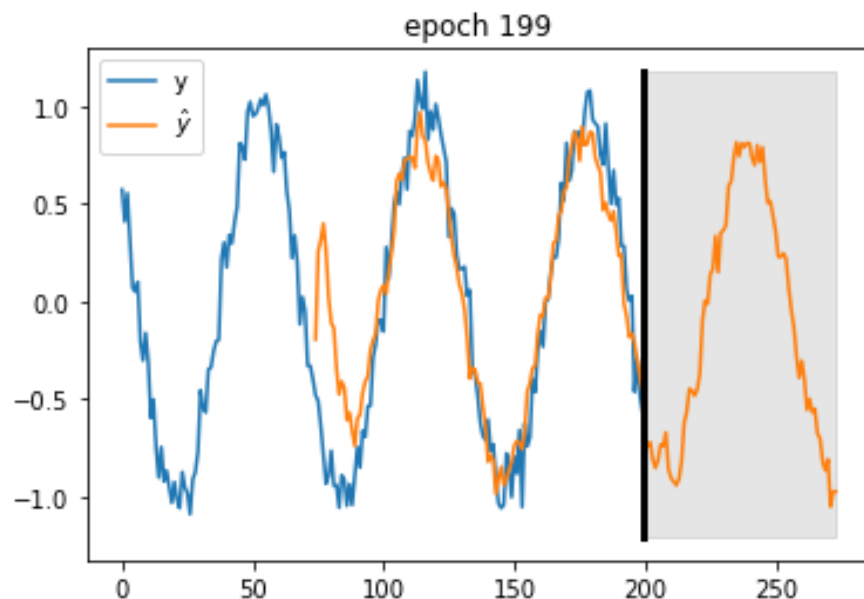




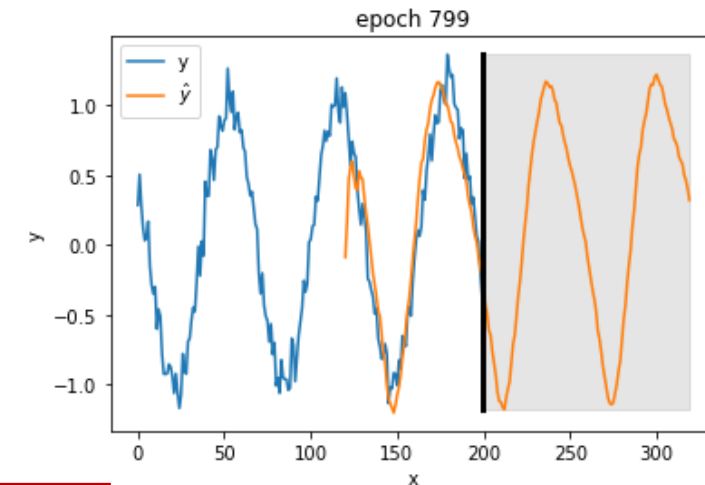
now: $Y(t) \rightarrow Y(t + dt)$

$dt = 75$

```
[lstm, dense1, dense2] = RunMyLSTM(Y_t, Y_t, n_neurons = 200,\n                                     n_epoch = 200, plot_each = 10, dt = dt,\n                                     momentum = 0.8, decay = 0.01,\n                                     learning_rate = 1e-3)
```



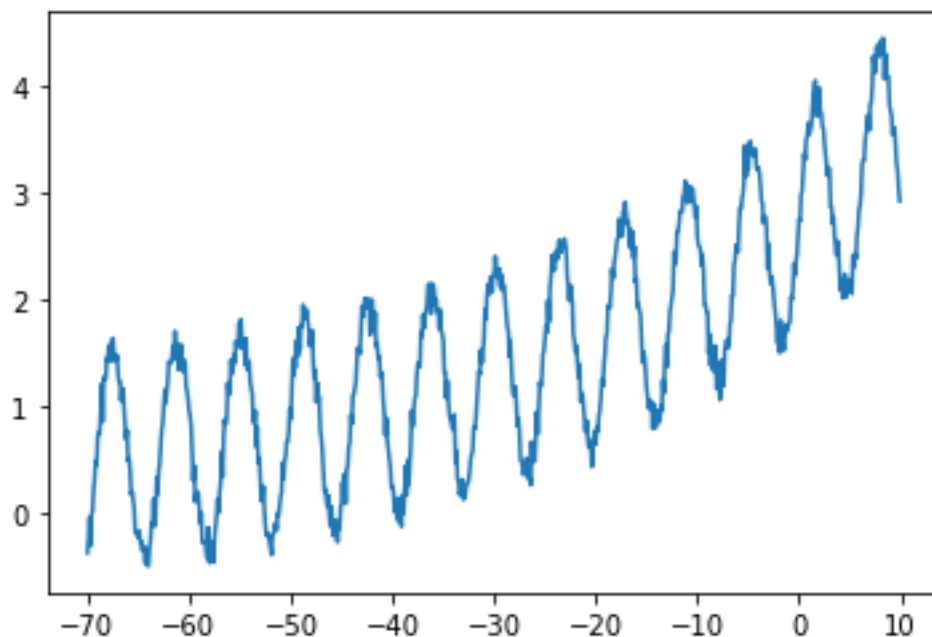
classical RNN



That is a lot better (and faster!) than the classical RNN!

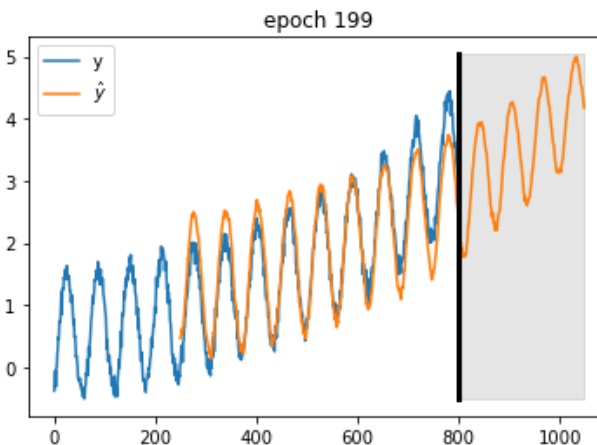
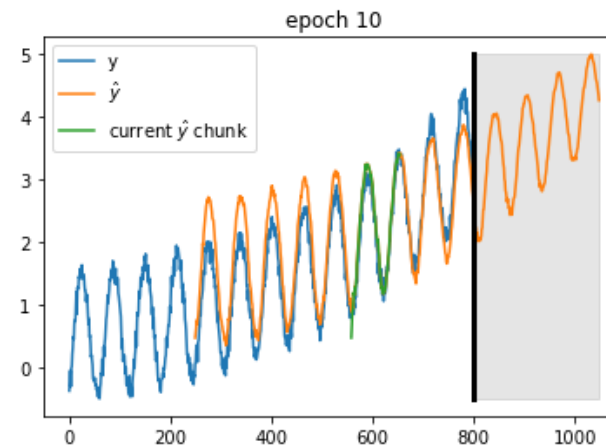
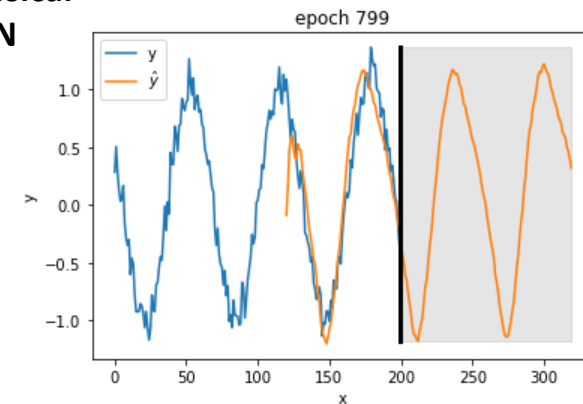
now: $Y(t) \rightarrow Y(t + dt)$

Let's challenge the LSTM and include long-term changes:


$$dt = 250$$

```
[lstm, dense1, dense2] = RunMyLSTM(Y_t, Y_t, n_neurons = 200,\n                                     n_epoch = 200, plot_each = 10, dt = dt,\n                                     momentum = 0.8, decay = 0.01,\n                                     learning_rate = 1e-5)
```

**classical
RNN**





now: $Y(t) \rightarrow Y(t + dt)$

Let's do the same thing with the **RNN** now:

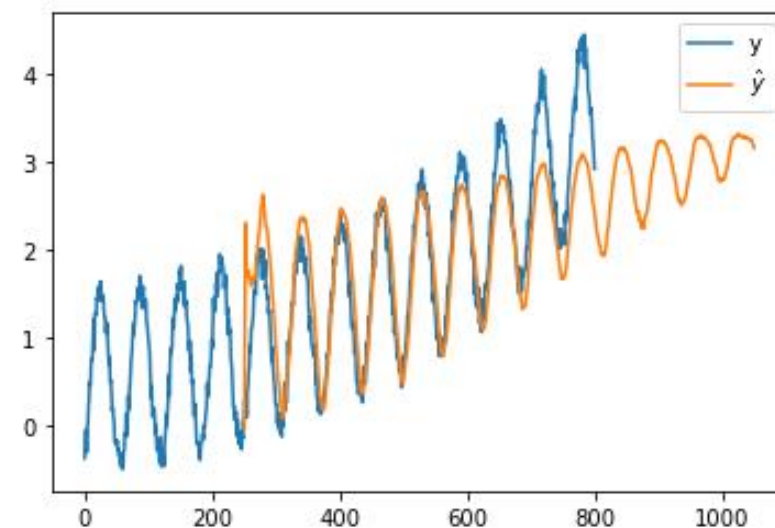
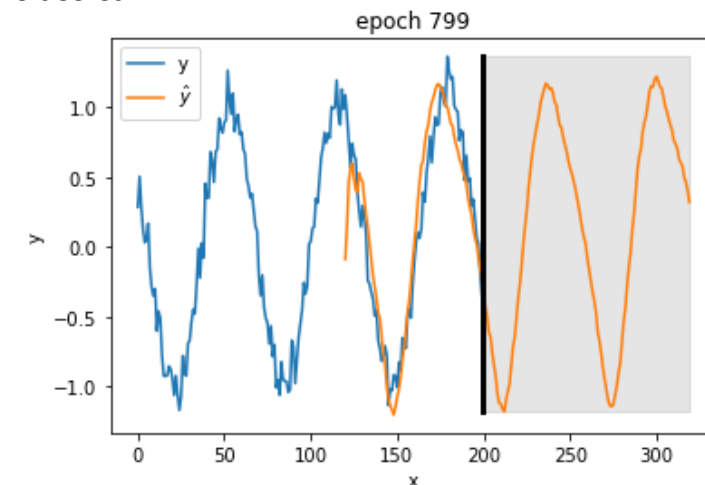
```
dt = 250
rnn = RunMyRNN(Y_t, Y_t, n_neurons = 200,\
                n_epoch = 200, plot_each = 10, dt = dt,\
                momentum = 0.8, decay = 0.01,\
                learning_rate = 1e-5)
```

```
Y_hat = ApplyMyRNN(Y_t, rnn)
```

```
X_plot = np.arange(0, len(Y_t))
X_plot_hat = np.arange(0, len(Y_hat)) + dt
```

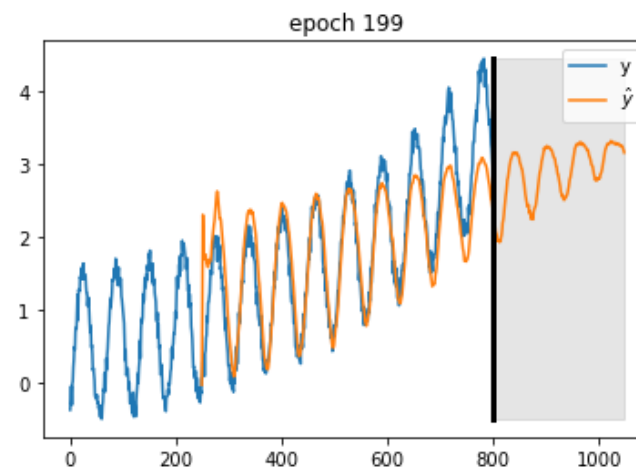
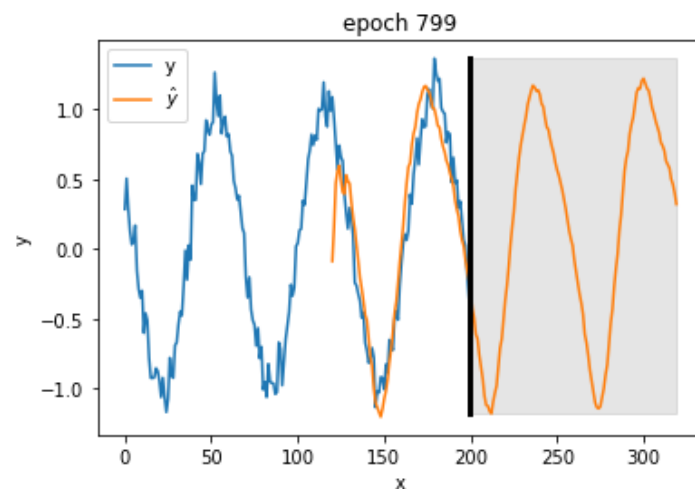
```
plt.plot(X_plot, Y_t)
plt.plot(X_plot_hat, Y_hat)
plt.legend(['y', '$\hat{y}$'])
plt.show()
```

classical RNN

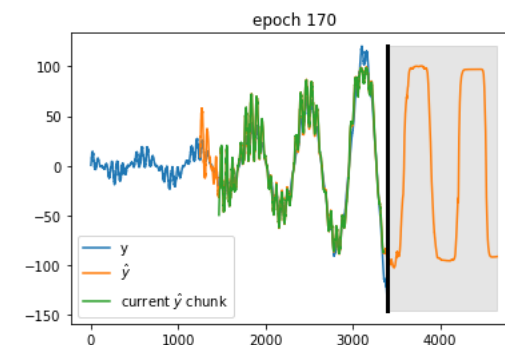
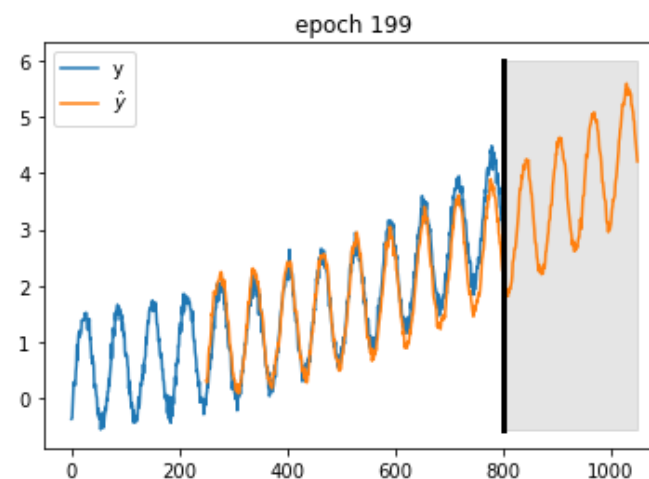
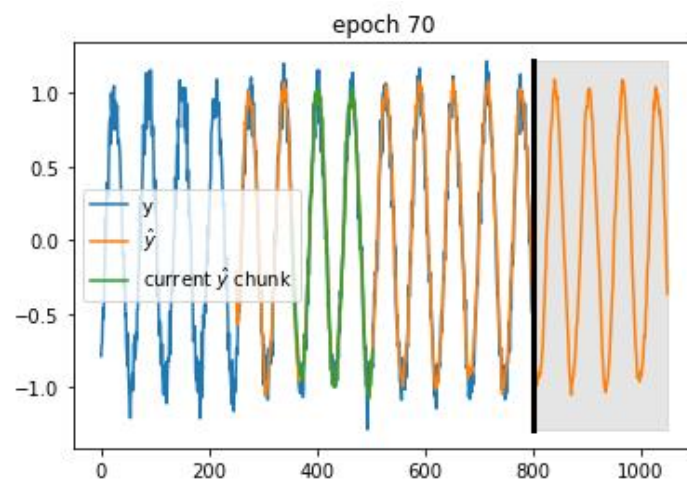




classical RNN



LSTM





Thank you very much for your attention!

