

# Unsupervised learning, Long Short Term Memory (LSTM), Transformer

Wen Yu

One object of Deep Learning: **extract features automatically**.  
We have no-label data  $x_1(k) \cdots x_n(k)$ , but want to find the latent (hidden features) inside  $x(k)$

$$(\lambda_1 \cdots \lambda_m) = PCA[x(k)]$$

But if we use

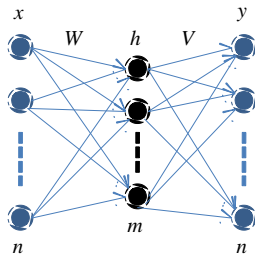
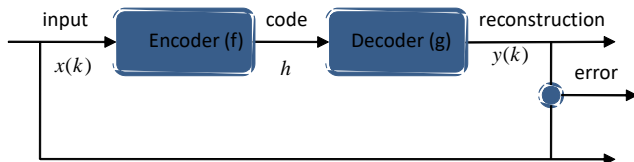
$$z = f(x), \quad \hat{x} = f^{-1}(z)$$

then

$$\hat{x} = f^{-1}[f(x(k))] = x(k)$$

So the maong  $F$  is regarded as a representation of the input  $x(k)$

# Autoencoder structure



# Autoencoder structure

- encoder: input layer  $\rightarrow$  hidden layer

$$z = f(x) = \phi_f(Wx)$$

- decoder: hidden layer  $\rightarrow$  output layer

$$\hat{x} = g(z) = \phi_g(Vz)$$

$\phi_f$  and  $\phi_g$  can be sigmoid functions, to simplify the computation (tied weights)

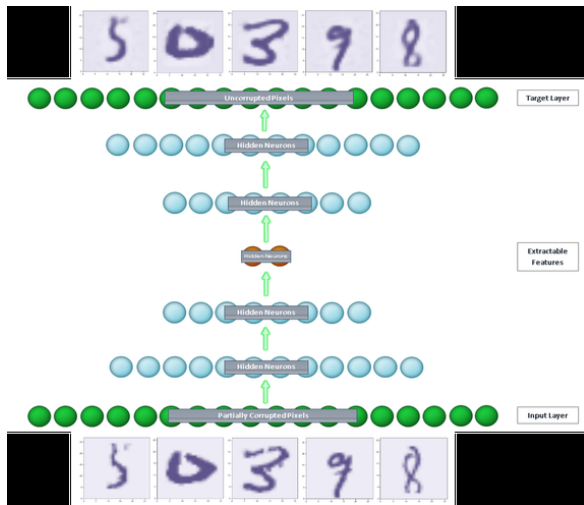
$$V = W^T$$

- decrease dimension

$$m < n$$

- When  $\phi_f$  is a linear function, autoencoder becomes PCA

# Unsupervised learning



# Autoencoder training

- The object: we hope the decoded variable  $y$  can approximate the original variable  $x$ , although it is encoded into  $h$ .
- Define the reconstruction errors as

squared error:  $e(k) = \|x(k) - \hat{x}(k)\|^2$

cross entropy:  $e(k) = \sum_{i=1}^n [x_i \log(\hat{x}_i) + (1 - x_i) \log(1 - \hat{x}_i)]$

- Loss functions

Average:  $L = \frac{1}{N} \sum_{k=1}^N |e(k)|, \quad \min_{\theta} L = \sum_{x \in S} e(x, g[f(c)])$

where the training data are  $x(1) \cdots x(N) \in S, \theta = [W, p, q]$

# Regularized autoencoder

Auto-encoder may perfectly construct the input without extracting any useful features. For example  $f(x) = x$ , or  $m \gg n$  (over-complete setting).

The loss functions should be modified.

L1 regularied (Lasso)

$$L = \sum_{x \in S} e(x, g[f(x)]) + \sum_{ij} |w_{ij}|$$

L2 regularied (Ridge)

$$L = \sum_{x \in S} e(x, g[f(x)]) + \lambda \sum_{ij} w_{ij}^2$$

where  $W = [w_{ij}]$ ,  $\lambda$  is the weigh decay factor

BP algorithm

$$\Delta W(k+1) = -\eta \frac{\partial L}{\partial W(k)}$$

can be used.

LS version

$$W(k) = \left( X^T X + \lambda I \right)^{-1} X^T Y$$

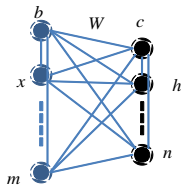


Geoffrey Hinton and Terry Sejnowski in 1985

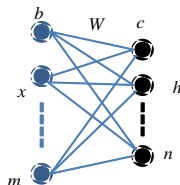
- stochastic
  - bi-direction network
- 1 Boltzmann machines with unconstrained connectivity have not proven useful for practical problems in machine learning or inference
  - 2 if the connectivity is properly constrained, the learning can be useful for practical problems

# Restricted Boltzmann Machines

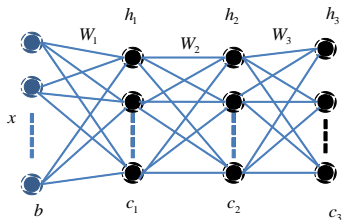
Undirected graphical model of the Restricted Boltzmann Machine (RBM)



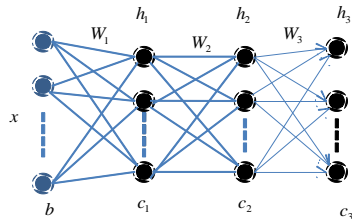
Boltzmann Machine



Restricted Boltzmann Machine



Deep Restricted Boltzmann Machines



Deep Belief Network

- no links between units of the same layer, only between input (or visible) units  $x_j$  and hidden units  $h_i$ . All hidden nodes are independent when the input nodes are known,

$$p(h | x) = \prod_{i=1}^n p(h_i | x)$$

all visible nodes are independent when the hidden nodes are known

$$p(x | h) = \prod_{j=1}^m p(x_j | h)$$

- all nodes have values of 0 or 1

- all  $x$  and  $h$  satisfy Boltzmann distribution (Gibbs distribution):

$$p(x) \propto e^{-\frac{E}{kT}}$$

where  $E$  is state energy, and  $k$  is the Boltzmann's constant,  $T$  is the thermodynamic temperature.

- The distribution shows that states with lower energy will always have a higher probability of being occupied than the states with higher energy (more random).
- The quantitative relationship between the probabilities of the two states being occupied

$$\frac{p_i}{p_j} = e^{-\frac{E_i - E_j}{kT}}$$

- Input  $(x_1 \cdots x_m)$  after RBM encoded become  $(h_1 \cdots h_n)$ ,  $i$ -th node in hidden layer, the probability of 1 is

$$p(h_i = 1 | x) = \phi [Wx + c]$$

where  $\phi$  is a sigmoid function

- How to calculate get  $h_i$ , if  $\phi [Wx + c] = a$ ,  $a \in (0, 1)$

$$h_i = \begin{cases} 1 & a < p(h_i = 1 | x) \\ 0 & a \geq p(h_i = 1 | x) \end{cases}$$

If  $\phi [Wx + c] = p(h_i = 1 | x) = 0.85$  means the probability of  $h_i = 1$  is 0.85.

Or when  $x$  is a random number between 0 and 1, it has 85% in  $(0, 0.85)$ . So when the random number is in  $(0, 0.85)$ , i.e.,  $a < p(h_i = 1 | x)$ , the event occurs,  $h_i = 1$

The energy of RBM is defined as

$$E(x, h) = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} x_j h_i - \sum_{j=1}^m b_j x_j - \sum_{i=1}^n c_i h_i$$

$w_{ij}$  are the weights associated with the connection between hidden unit  $h_i$  and visible unit  $x_j$

$b_j$  are the bias weights for the visible units  $x_j$

$c_i$  are the bias weights for the hidden units  $h_i$

The energy depends on each node and connection

The joint probability of visual node and the hidden node is

$$p(x, h) = \frac{e^{-E(x, h)}}{Z} \leq 1$$

where the normalizing factor  $Z$  is called the partition function

$$Z = \sum_{x, h} e^{-E(x, h)}$$

It is Boltzmann distribution.

The distance (difference) from the probability distribution of RBM  $p(x)$ , to the probability distribution of the input  $q(x)$  is defined as

$$KL(p, q) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$

WHERE  $KL(p, q)$  is the the Kullback-Liebler divergence, so

$$KL(p, q) = \sum_x q(x) \log q(x) - \sum_x q(x) \log p(x)$$

the first term  $\sum_x q(x) \log q(x)$  is entropy of  $x$ ,



# Training error

The second term cannot be obtained directly, should be estimated by Monte Carlo method, or estimated by

$$\sum_x q(x) \log p(x) \approx \frac{1}{l} \sum_{x \in X(l)} \log p(x)$$

where  $l$  is selected samples.

Now

$$\min KL(p, q) = \min \left[ \sum_x q(x) \log q(x) - \sum_x q(x) \log p(x) \right] \rightarrow \max \sum_x \log p$$

If RBM have some random states  $(x, h)$ , the probability in visual node arrives maximum from  $h$  to  $x$  process, i.e., the error arrive minimum

$$\frac{\partial L(\theta)}{\partial \theta} = \sum_v \frac{\partial \log p(v)}{\partial \theta}$$

$$\theta(k+1) = \theta(k) - \eta \frac{\partial [-\log p(\mathbf{x})]}{\partial \theta(k)}$$

where  $\theta = [W, b, c]$

$$\begin{aligned}\frac{\partial L(\theta)}{\partial w_{ij}} &= \sum_{\mathbf{v}} p(h_i = 1 | \mathbf{v}) v_j - \sum_{\mathbf{v}} \frac{1}{l} \sum_{k=1}^l p(h_i = 1 | \mathbf{v}_{y_k}) v_{y_k, j} \\ \frac{\partial L(\theta)}{\partial b_j} &= \sum_{\mathbf{v}} v_j - \sum_{\mathbf{v}} \frac{1}{l} \sum_{k=1}^l v_{y_k, j} \\ \frac{\partial L(\theta)}{\partial c_i} &= \sum_{\mathbf{v}} p(h_i = 1 | \mathbf{v}) - \sum_{\mathbf{v}} \frac{1}{l} \sum_{k=1}^l p(h_i = 1 | \mathbf{v}_{y_k})\end{aligned}$$

- $\sum_{\mathbf{v}} v_j$  and  $\sum_{\mathbf{v}} p(h_i = 1 | \mathbf{v})$  mean to summarize all input samples
- $\sum_{\mathbf{v}} \frac{1}{l} \sum_{k=1}^l v_{y_k, j}$  means to summarize  $l$  output samples

# LSTM (Long Short Term Memories)

## Long-Term Dependencies

- RNNs can learn to use the past information, where the gap between the relevant information and the place are small.
- In theory, RNNs are absolutely capable of handling such “long-term dependencies” by picking parameters.
- In practice, RNNs don't seem to be able to learn them.
- The vanishing gradient problem for RNNs

# Long Short Term Memory networks

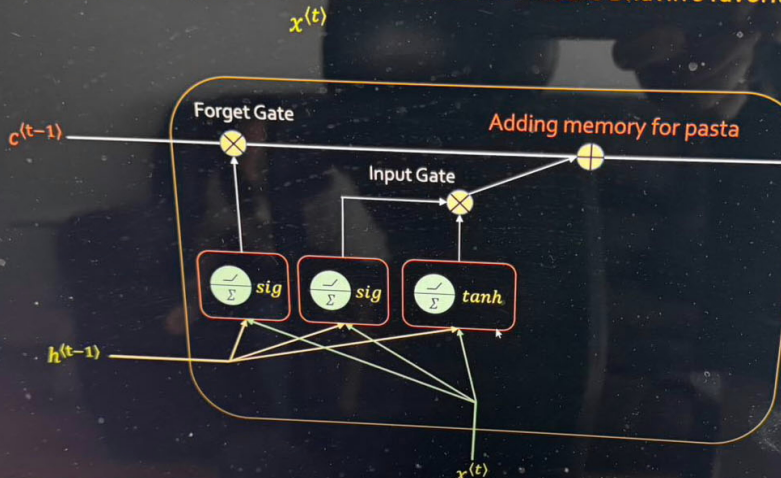
LSTM (Long Short Term Memories) is a special type of recurrent neural network structure, its hidden layer has a special structure. By introducing a gating mechanism to solve the gradient vanishing problem of RNN, it can learn long-distance dependencies.

In theory, RNNs can indeed connect long-term dependencies and solve such problems. But unfortunately in practice, RNNs cannot solve this problem. Hochreiter (1991) and Bengio, et al. (1994) have studied this problem in depth and found that RNNs are indeed difficult to solve this problem.

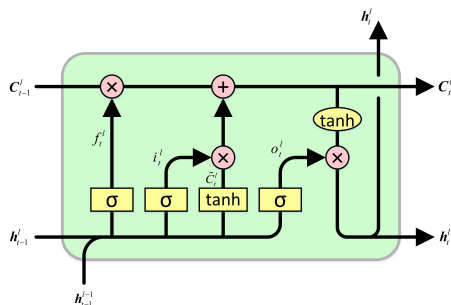
*S. Hochreiter and J. Schmidhuber. Long short-term memory. Neural computation, Vol.9, No.8, 1997*

# LSTM

Dhaval eats samosa almost everyday, it shouldn't be hard to guess that his favorite  
Bhavin however is a lover of pasta and cheese that means Bhavin's favorite



Finally



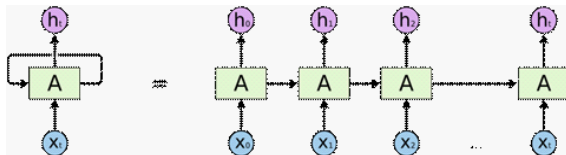
$$y_k = \sigma(W_{O,k} [y_{k-1}, u_k]) \tanh(x_k)$$

where

$$x_k = \sigma(W_{F,k} [y_{k-1}, u_k]) x_{k-1} + \sigma(W_{I,k} [y_{k-1}, u_k]) \tanh(W_{X,k} [y_{k-1}, u_k])$$

# Recurrent neural network (RNN)

- LSTM is a recurrent neural network (RNN)
- A RNN can be thought of as multiple copies of the same network, each passing a message to a successor
- RNNs are related to sequences and lists. The architecture of RNN is suitable for such data.

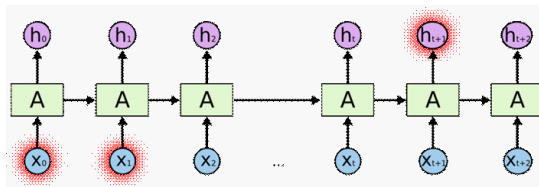
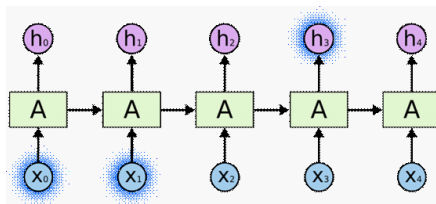




# RNN

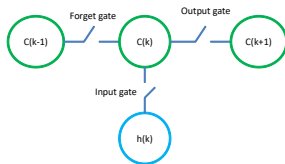
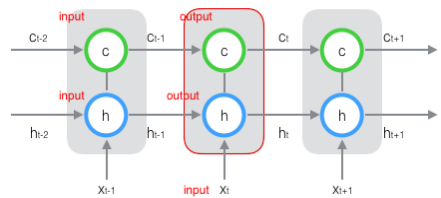
Connecting previous information to the present task depends many factors:  
gap length.

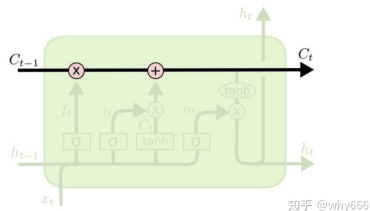
As the gap grows, RNNs become unable to learn to connect the  
information



# LSTM-structure

- 1) Cell state  $c$  (keep LongTerm Memory), while RNN only has hidden state  $h$  (keep Short Term Memory)
- 2) Gate is a way to let information through. LSTM has three gates to protect and control the cell state: forget gate, input gate, output gate



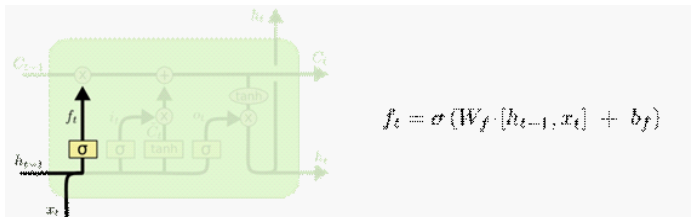


The transfer of cell state is like a conveyor belt, with vectors passing through the entire cell, with only a few linear operations. This structure makes it easy to pass information through the entire cell without changing it. (Achieving long-term memory retention)

# Forget gate

It looks at  $h_{k-1}, x_k$ , and outputs a number between 0 and 1 for each number in the cell state  $C_{t-1}$ .

1 represents “completely keep this” while 0 represents “completely get rid of this.”



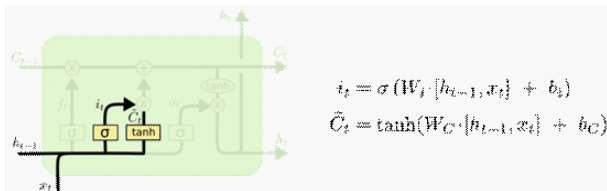
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$f_t = \sigma(W_f [h_{t-1}, x_t]),$$

# Input gate

For information store, to decide what new information we're going to store in the cell state. This has two parts

- 1 a sigmoid layer called the “input gate layer” decides which values we'll update.
- 2 a tanh layer creates a vector of new candidate values,  $\tilde{C}_t$ , that could be added to the state.



$$i_t = \sigma(W_i [h_{t-1}, x_t] + b_i)$$

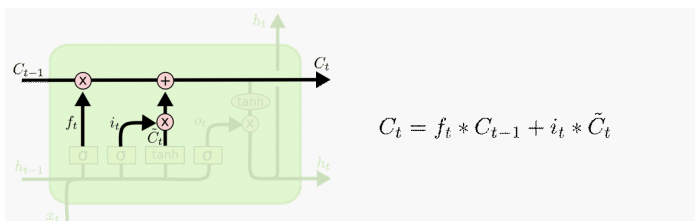
$$\tilde{C}_t = \tanh(W_c [h_{t-1}, x_t] + b_c)$$

$$i_k = \sigma(W_i [h_{t-1}, x_t])$$

$$\tilde{C}_t = \tanh(W_c [h_{t-1}, x_t])$$

# Cell update

We multiply the old state  $C_{t-1}$  by  $f_t$ , forgetting the things we decided to forget earlier. Then we add  $i_t \tilde{C}_t$ .



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

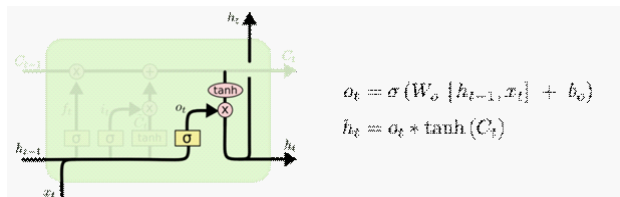
$$C_t = f_t C_{t-1} + i_t \tilde{C}_t,$$

# Output gate

This output will be based on our cell state, but will be a filtered version.

- 1 we run a sigmoid layer which decides what parts of the cell state we're going to output.
- 2 we put the cell state through tanh (to push the values to be between -1 and 1 ) and multiply it by the output of the sigmoid gate,

We only output the parts we decided to.



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$\tilde{h}_t = o_t * \tanh(C_t)$$

$$o_t = \sigma(W_o [h_{t-1}, x_t])$$

$$h_t = o_t \tanh(C_t)$$

$$\begin{aligned}
 f_t &= \sigma(W_f [h_{t-1}, x_t]) \\
 i_t &= \sigma(W_i [h_{t-1}, x_t]) \\
 \tilde{C}_t &= \tanh(W_c [h_{t-1}, x_t]) \\
 C_t &= f_t C_{t-1} + i_t \tilde{C}_t \\
 o_t &= \sigma(W_o [h_{t-1}, x_t]) \\
 h_t &= o_t \tanh(C_t)
 \end{aligned}$$

So

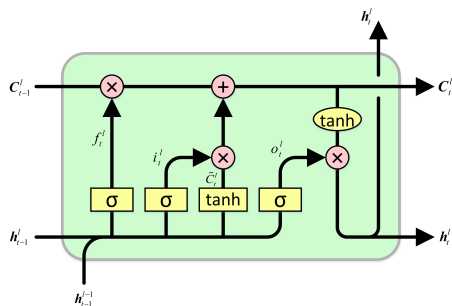
$$h_t = \sigma(W_o [h_{t-1}, x_t]) \tanh(C_t)$$

where

$$\begin{aligned}
 C_t &= f_t C_{t-1} + i_t \tilde{C}_t \\
 &= \sigma(W_f [h_{t-1}, x_t]) C_{t-1} + \sigma(W_i [h_{t-1}, x_t]) \tanh(W_c [h_{t-1}, x_t])
 \end{aligned}$$



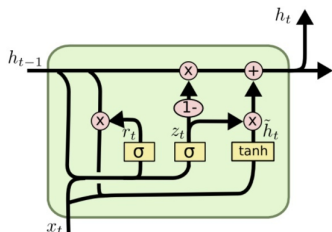
Finally



$$y_k = \sigma(W_{O,k} [y_{k-1}, u_k]) \tanh(x_k)$$

where

$$x_k = \sigma(W_{F,k} [y_{k-1}, u_k]) x_{k-1} + \sigma(W_{I,k} [y_{k-1}, u_k]) \tanh(W_{X,k} [y_{k-1}, u_k])$$



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

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

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

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Gated Recurrent Unit (GRU), Cho, et al. (2014).

It combines the forget gate and the input gate into a single “update gate.”

1) Forget gate:

$$F_k = \sigma (W_{F,k} [y_{k-1}, u_k])$$

2) Input gate:

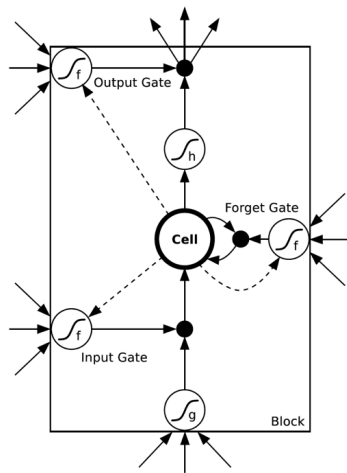
$$\begin{aligned} I_k &= \sigma (W_{I,k} [y_{k-1}, u_k]_I) \\ \tilde{x}_k &= \tanh (W_{X,k} [y_{k-1}, u_k]) \\ x_k &= F_k x_{k-1} + I_k \tilde{x}_k \end{aligned}$$

So

$$y_k = y_{k-1} + \sigma (W_{F,k} [y_{k-1}, u_k]) (\tanh (W_{X,k} [I_k y_{k-1}]) - y_{k-1})$$

It merges the cell state and hidden state. GRU is simpler than standard LSTM models.

# LSTM



Use coupled forget and input gates. We only forget when we're going to input something in its place.

We only input new values to the state when we forget something older.

$$x_k = F_k x_{k-1} + (1 - F_k) \tilde{x}_k$$

Depth Gated RNNs by Yao, et al. (2015)

Clockwork RNNs by Koutnik, et al. (2014).

- Greff, et al. (2015) do a nice comparison of popular variants, finding that they're all about the same.
- Jozefowicz, et al. (2015) tested more than ten thousand RNN architectures, finding some that worked better than LSTMs



"peephole connections: the gate layers receive **cell state**

$$F_k = \sigma(W_{F,k} [x_{k-1}, y_{k-1}, u_k] + b_f)$$

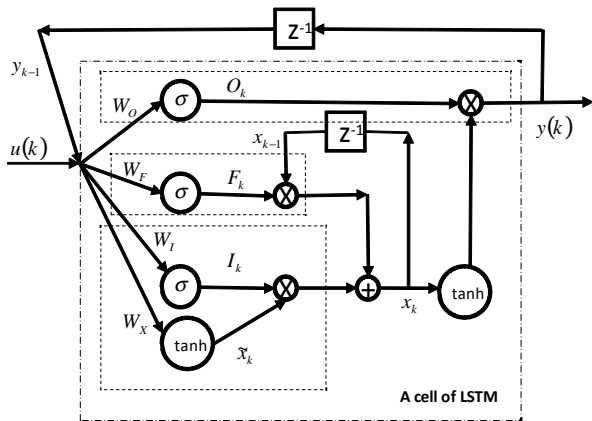
$$I_k = \sigma(W_{I,k} [x_{k-1}, y_{k-1}, u_k] + b_I)$$

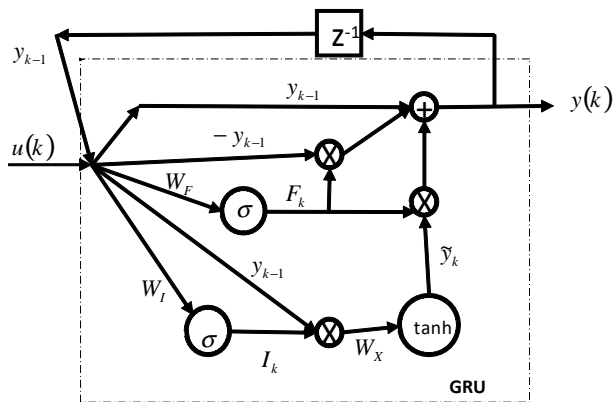
$$O_k = \sigma(W_{O,k} [x_{k-1}, y_{k-1}, u_k] + b_O)$$

$$\tilde{x}_k = \tanh(W_{X,k} [y_{k-1}, u_k] + b_x)$$



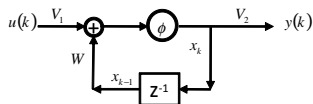
# LSTM in the form of system



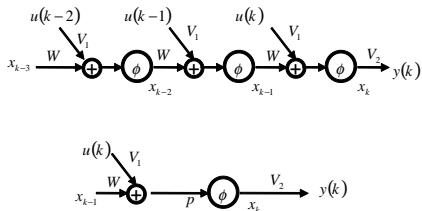
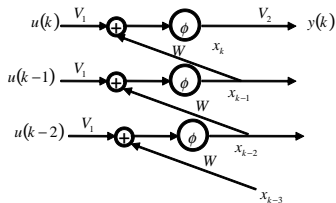


# Training of RNN

Recurrent NN  $\rightarrow$  feedforward NN



Unfold



Recurrent cell

$$\begin{aligned}y(k) &= V_2(k) x(k) \\x(k) &= \phi [W(k) x(k-1) + V_1(k) u(k)]\end{aligned}$$

It can be unfolded into  $m$  steps

$\vdots$

$$\begin{aligned}y(k-m) &= V_2(k) x(k-m) \\x(k-m) &= \phi [W(k) x(k-m-1) + V_m(k) u(k-m)]\end{aligned}$$

## Backpropagation Through Time (BPTT)

$y^*(k)$  is the desired output, the instantaneous of the squared output errors is

$$J(k) = \frac{1}{2} [y(k) - y^*(k)]^2 = e_o^2(k)$$

The gradient decent is

$$w(k+1) = w(k) - \eta \frac{\partial J(k)}{\partial w(k)}$$
$$\frac{\partial J}{\partial w} = \sum_{k=1}^m \frac{\partial J}{\partial w(k)}$$

# Training of output layer

So

$$V_2(k+1) = V_2(k) - \eta \frac{\partial J(k)}{\partial V_2(k)}$$

Because  $y(k) = V_2(k)x(k)$

$$\begin{aligned} \frac{\partial J(k)}{\partial V_2(k)} &= \frac{\partial J(k)}{\partial e_o(k)} \frac{\partial e_o(k)}{\partial y(k)} \frac{\partial y(k)}{\partial V_2(k)} \\ &= e_o(k) \times 1 \times x(k) \end{aligned}$$

$$V_2(k+1) = V_2(k) - \eta e_o x(k)$$

So

$$V_1(k+1) = V_1(k) - \eta \frac{\partial J(k)}{\partial V_1(k)}$$

Because  $y(k) = V_2(k)x(k)$ ,  $p = V_1(k)u(k) + Wx(k-1)$

$$\begin{aligned} \frac{\partial J(k)}{\partial V_1(k)} &= \frac{\partial J(k)}{\partial e_o(k)} \frac{\partial e_o(k)}{\partial y(k)} \frac{\partial y(k)}{\partial x(k)} \frac{\partial x(k)}{\partial p(k)} \frac{\partial p(k)}{\partial V_1(k)} \\ &= e_o(k) \times 1 \times V_2(k) \times \phi' \times u(k) \end{aligned}$$

$$V_1(k+1) = V_1(k) - \eta e_o(k) V_2(k) \phi' u(k)$$

$$\begin{aligned}\frac{\partial J(k)}{\partial W} &= \frac{\partial J(k)}{\partial e_o(k)} \frac{\partial e_o(k)}{\partial y(k)} \frac{\partial y(k)}{\partial x(k)} \frac{\partial x(k)}{\partial p(k)} \frac{\partial p(k)}{\partial W} \\ &= e_o(k) \times 1 \times V_2(k) \times \phi' \times x(k-1)\end{aligned}$$

$$W(k+1) = W(k) - \eta e_o(k) V_2(k) \phi' x(k-1)$$



Similar

$$e_1(k) = e_o(k) V_2(k) \phi'$$

$$V_1(k+1) = V_1(k) - \eta e_1(k) u(k)$$

$$W(k+1) = W(k) - \eta e_1(k) x(k-1)$$

$$e_2(k) = e_1(k) W(k) \phi'$$

$$V_1(k+1) = V_1(k) - \eta e_2(k) u(k-1)$$

$$W(k+1) = W(k) - \eta e_2(k) x(k-2)$$

# Backpropagation Through Time (BPTT)

Finally

$$V_1(k+1) = V_1(k) - \eta e_1(k) u(k)$$

$$V_1(k) = V_1(k-1) - \eta e_2(k) u(k-1)$$

$$V_1(k-1) = V_1(k-2) - \eta e_3(k) u(k-2)$$

$$e_1(k) = e_o(k) V_2(k) \phi'$$

$$e_2(k) = e_1(k) W(k) \phi'$$

$$e_3(k) = e_2(k) W(k) \phi'$$

$$W(k+1) = W(k) - \eta e_1(k) x(k-1)$$

$$W(k) = W(k-1) - \eta e_2(k) x(k-2)$$

$$W(k-1) = W(k-2) - \eta e_3(k) x(k-3)$$

# Backpropagation Through Time (BPTT)

So

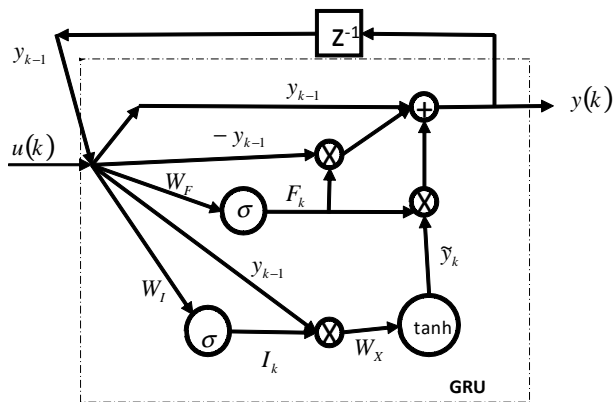
$$\begin{aligned}W(k+1) &= W(k) - \eta \sum_{i=1}^M e_i(k) x(k-i) \\V_1(k+1) &= V_1(k) - \eta \sum_{i=1}^M e_i(k) u(k-i+1) \\V_2(k+1) &= V_2(k) - \eta e_o x(k)\end{aligned}$$

where

$$e_i(k) = e_{i-1}(k) W(k) \phi'$$

If  $\|W(k) \phi'\| > 1$ , it is gradient blow up. If  $\|W(k) \phi'\| < 1$ , it is gradient vanish. We need to add a constraint as

when  $\|W(k) \phi'\| \approx 1$ , the layer is activated



Recurrent part

$$y_k = y_{k-1} - F_k y_{k-1} + F_k \tilde{y}_k$$

Feedforward part

$$\begin{aligned} F_k &= \sigma(W_{F,k} [y_{k-1}, u_k]) \\ \tilde{y}_k &= \phi(W_{X,k} [I_k y_{k-1}, u_k]) \\ I_k &= \sigma(W_{I,k} [y_{k-1}, u_k]) \end{aligned}$$

In state space

$$\begin{aligned} x_k &= x_{k-1} - \sigma(W_{F,k} [x_{k-1}, u_k]) x_{k-1} \\ &\quad + \sigma(W_{F,k} [x_{k-1}, u_k]) \phi(W_{X,k} [\sigma(W_{I,k} [x_{k-1}, u_k]) x_{k-1}]) \end{aligned}$$

$$J(k) = \frac{1}{2} [y(k) - d(k)]^2 = e_i^2(k)$$
$$W(k+1) = W(k) - \eta \frac{\partial J(k)}{\partial W(k)}$$

If

$$y(k) = \phi [W(k) y(k-1) + V_1(k) u(k)]$$

BPTT is

$$W(k+1) = W(k) - \eta \sum_{i=1}^M e_i(k) y(k-i)$$
$$V_1(k+1) = V_1(k) - \eta \sum_{i=1}^M e_i(k) u(k-i+1)$$
$$e_i(k) = e_{i-1}(k) W(k) \phi'$$

For  $W_{F,k}$ , because

$$y_k = y_{k-1} - \sigma(W_{F,k} [y_{k-1}, u_k]) y_{k-1} + F_k \tilde{y}_k$$

BPTT

$$W_F(k+1) = W_F(k) - \eta \sum_{i=1}^M e_F(k) [y_{k-i}, u_{k-i}]$$
$$e_F(k) = e_{F-1}(k) W_F(k) \sigma' y_{k-i}$$

or

$$W_F(k+1) = W_F(k) - \eta \sigma' \sum_{i=1}^M e_{F-1}(k) W_F(k) [y_{k-i}, u_{k-i}]$$

For  $W_X$  and  $W_I$ ,

$$y_k = y_{k-1} - F_k y_{k-1} + F_k \phi (W_{X,k} [\sigma (W_{I,k} [x_{k-1}, u_k]) x_{k-1}])$$

BPTT

$$\begin{aligned}W_X(k+1) &= W_X(k) - \eta \sum_{i=1}^M e_X(k) [I_k y_{k-1}, u_k] \\e_X(k) &= e_{X-1}(k) W_X(k) \phi' \sigma (W_{F,k} [y_{k-1}, u_k]) \\W_I(k+1) &= W_I(k) - \eta \sum_{i=1}^M e_I(k) y_{k-1} \\e_I(k) &= e_{I-1}(k) W_X(k) \phi' \sigma' \sigma W_I(W_{F,k} [y_{k-1}, u_k])\end{aligned}$$



# LSTM for nonlinear system identification

$$y(k) = \Phi[y(k-1), \dots, y(k-n_y), u(k), \dots, u(k-n_u)] \quad (1)$$

Simulation model,

$$\hat{y}(k) = F[u(k), \dots, u(k-n)] \quad (2)$$

$$\hat{y}(k) = N[\hat{y}(k-1), \dots, \hat{y}(k-m), u(k), \dots, u(k-n)] \quad (3)$$

It is the parallel identification model.

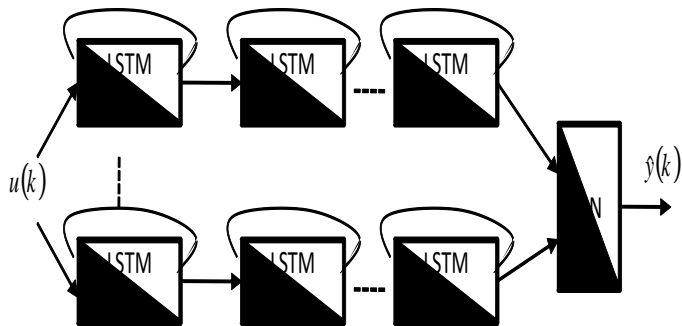
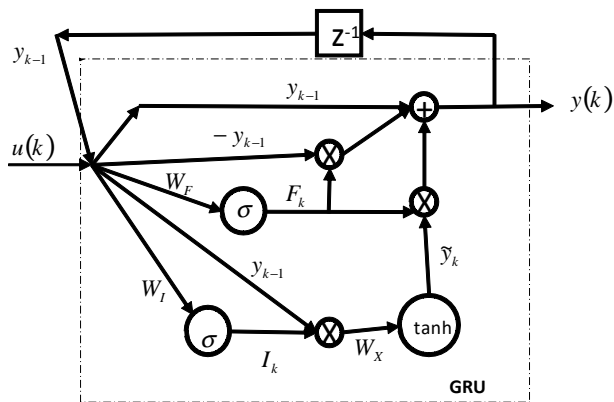


Figure: LSTM for nonlinear system modeling



$$e(k) = \hat{y}(k) - y(k)$$

The dynamic of the GRU is

$$x_{k+1} = x_k - \sigma(W_1 x_k + W_u u_k) x_k + \sigma(W_1 x_k + W_u u_k) \phi(W_2 [\sigma(W_3 x_k) x_k])$$

The GRU can be transformed into two parts: RNN and FNN (feedforward NN)

$$\begin{aligned} \text{RNN: } & x(k+1) = Ax(k) - \sigma[W_1 x(k) + W_u u(k)] \tilde{u}(k) \\ \text{FNN: } & \tilde{u}(k) = \phi(W_2 \sigma[W_3 x(k)] x(k)) + x(k) \end{aligned} \quad (4)$$

## Theorem

For the feedforward part of GRU, FNN, the following backpropagation-like algorithm can make identification error  $e_F(k)$  bounded

$$\begin{aligned} W_{2,k+1} &= W_{2,k} - \eta_k e_F(k) \phi' x^2(k) \sigma \\ W_{3,k+1} &= W_{3,k} - \eta_k e_F(k) \phi' \sigma' W_{2,k} x^2(k) \end{aligned} \quad (5)$$

where  $\eta_k = \frac{\eta}{1 + \|\phi' x^2(k) \sigma\|^2 + \|\phi' \sigma' W_{2,k} x^2(k)\|^2}$ . The average of the identification error satisfies

$$J = \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{k=1}^T e_F^2(k) \leq \frac{\eta}{\pi} \bar{\delta}_1 \quad (6)$$

where  $\pi = \frac{\eta}{1 + \kappa} \left[ 1 - \frac{\kappa}{1 + \kappa} \right] > 0$ ,

$\kappa = \max \left( \|\phi' x^2(k) \sigma\|^2 + \|\phi' \sigma' W_{2,k} x^2(k)\|^2 \right)$ ,  $\bar{\delta}_1 = \max [\delta_1^2(k)]$ ,

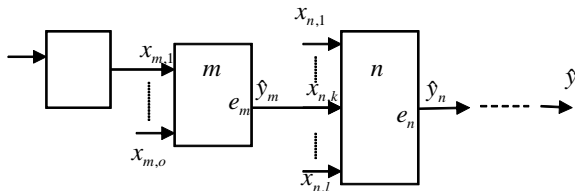


Figure: Error backpropagation

The performance index is defined as

$$J = \frac{1}{2} e_o^2 \quad e_o = \hat{y} - y, \quad \hat{y} = \phi [W_{N,k} \hat{y}_i]$$

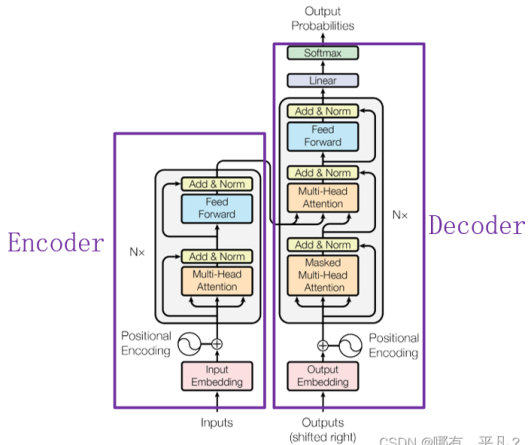
$$e_m = \frac{\partial \sigma}{\partial t} w_u e_n \quad (7)$$

the training law is the gradient descent

$$W_N(k+1) = W_N(k) - \eta \hat{y}_{i,q}(k) e_o(k) \quad (8)$$

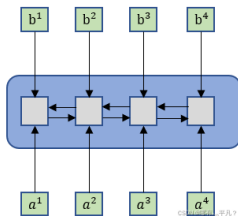
# Positional Encoding

Seq2seq



CSDN @哪有...平凡?

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin Attention Is All You Need, *Neural Information Processing Systems*, 2017



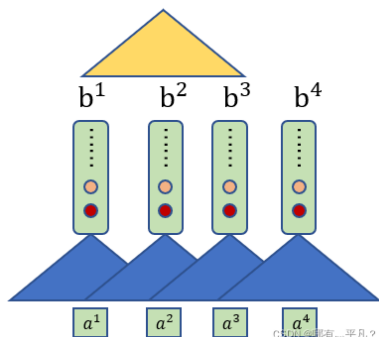
The problem of RNN

Input series  $\{a_1, a_2, a_3, a_4\}$ , the output series  $\{b_1, b_2, b_3, b_4\}$ . If we want  $b_3$  we need to input  $\{a_1, a_2, a_3\}$ , we cannot obtain all  $b_1, b_2, b_3, b_4$  at the same time.



# Parallel

One solution is to us filters (CNN)



Parallel serieses  $\{b_1, b_2, b_3, b_4\}$  with  $\{a_1, a_2\}$ ,  $\{a_1, a_2, a_3\}$ ,  $\{a_2, a_3, a_4\}$ ,  $\{a_3, a_4\}$

But we lost long term memory

# Self attention

Each input  $x_i$  is multiplied by  $W$ , we have  $(q, k, v)$

$$q_i = W^q a_i, \quad k_i = W^k a_i, \quad v_i = W^v a_i$$

where

$$a_i = Wx_i$$

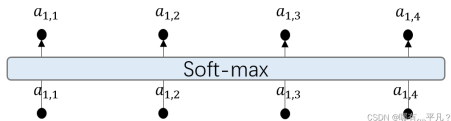
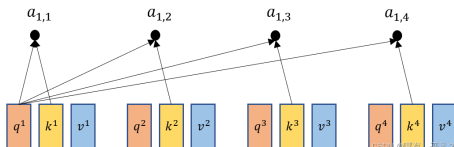
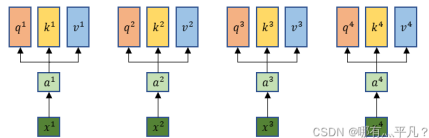
then we calculate dot product

$$a_{1,i} = \frac{1}{\sqrt{d}} q_1 k_i$$

then we apply soft-max

$$\hat{a}_{1,i} = \frac{e^{a_{1,i}}}{\sum_j e^{a_{1,j}}}$$

# Self attention

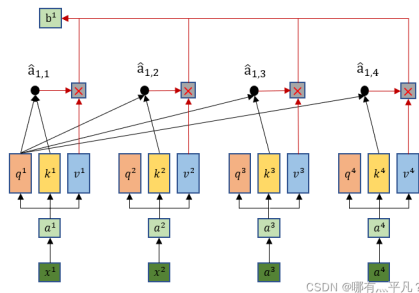


# Self attention

$\hat{a}_{1,i}$  is multiplied by  $v_i$

$$b_1 = \sum_i \hat{a}_{1,i} v_i$$

so  $b_1$  have all information of  $x_i$



Also: if  $b_1$  only pay attention on  $\{x_1, x_2\}$ , then we can set  $\hat{a}_{1,3} = \hat{a}_{1,4} = 0$

# Parallel

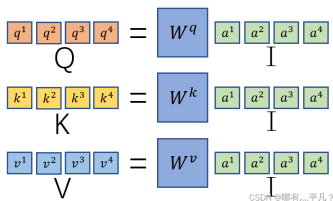
If we use vector form for,

$$q_i = W^q a_i, \quad k_i = W^k a_i, \quad v_i = W^v a_i, \quad a_i = W x_i$$

$$Q = [q_1, \dots, q_4], \quad K = [k_1, \dots, k_4], \quad V = [v_1, \dots, v_4]$$

The input is  $X$ , the output

$$A = XW, \quad Q = W^q A, \quad K = W^k A, \quad V = W^v A$$



$$\hat{A} = KQ, \quad \hat{A} = \text{soft max}(\hat{A})$$

$$\begin{array}{|c|c|c|c|} \hline a_{1,1} & a_{2,1} & a_{3,1} & a_{4,1} \\ \hline a_{1,2} & a_{2,2} & a_{3,2} & a_{4,2} \\ \hline a_{1,3} & a_{2,3} & a_{3,3} & a_{4,3} \\ \hline a_{1,4} & a_{2,4} & a_{3,4} & a_{4,4} \\ \hline \end{array} = \begin{array}{|c|} \hline k^1 \\ \hline k^2 \\ \hline k^3 \\ \hline k^4 \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline q^1 & q^2 & q^3 & q^4 \\ \hline \end{array}$$

A
K
Q

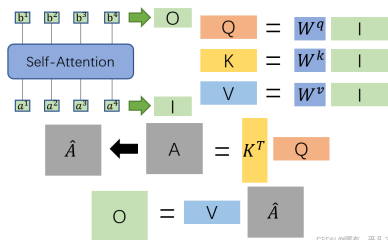
$$B = V\hat{A},$$

$$\begin{array}{|c|c|c|c|} \hline b^1 & b^2 & b^3 & b^4 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline v^1 & v^2 & v^3 & v^4 \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \hat{a}_{1,1} & \hat{a}_{2,1} & \hat{a}_{3,1} & \hat{a}_{4,1} \\ \hline \hat{a}_{1,2} & \hat{a}_{2,2} & \hat{a}_{3,2} & \hat{a}_{4,2} \\ \hline \hat{a}_{1,3} & \hat{a}_{2,3} & \hat{a}_{3,3} & \hat{a}_{4,3} \\ \hline \hat{a}_{1,4} & \hat{a}_{2,4} & \hat{a}_{3,4} & \hat{a}_{4,4} \\ \hline \end{array}$$

O
V
 $\hat{A}$

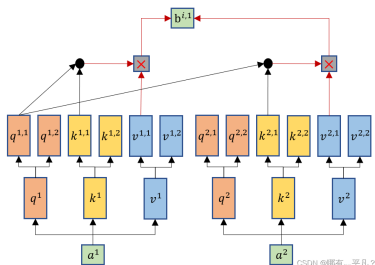
So

$$O = V\hat{A},$$



# Multi-head Self-attention

$(q, k, v)$  is divided into several, then each of them are applied to self attention





# Positional Encoding

Positions information: each element should has its own position information  $e_i$

