

Diagram illustrating a feedforward neural network structure with three layers:

- Input Layer:** Nodes  $x_1$  and  $x_2$ .
- Hidden Layer:** Nodes  $y_1$  and  $y_2$ .
- Output Layer:** Node  $y_n$ .

Connections (Weights) are shown between layers:

- From  $x_1$  to  $y_1$  (weight  $w_{11}$ )
- From  $x_1$  to  $y_2$  (weight  $w_{12}$ )
- From  $x_2$  to  $y_1$  (weight  $w_{21}$ )
- From  $x_2$  to  $y_2$  (weight  $w_{22}$ )
- From  $y_1$  to  $y_n$  (weight  $w_{31}$ )
- From  $y_2$  to  $y_n$  (weight  $w_{32}$ )

سین دایم به زبان سینی:

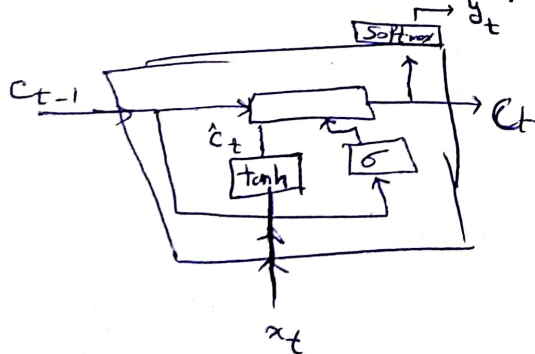
$$y_t = g_2(w_{y^a} a_t + b_y)$$

The diagram illustrates the internal structure of a Gated Recurrent Unit (GRU). It shows an input  $x_t$  and a hidden state  $h_{t-1}$  entering a 'tanh' block. The output of the 'tanh' block is multiplied by  $h_{t-1}$  to produce the new hidden state  $h_t$ . A 'sigmoid' block also receives  $x_t$  and  $h_{t-1}$  as input, and its output is multiplied by  $h_t$  to produce the final hidden state  $h_t$ .

✓ حال برای حل این مشکل ~~کمی~~ ماه GRU استفاده شود که حافظه بیشتری دارد.

حال برای سینه دادن state می دهی و می جاری در state میانی را در هم concept می کنیم و در وزن  $W$  ضرب می کنیم  
(به انتی های جبهه  $e$  می گوئیم که مختلف cell است)

و یک سیگنال اضافه شد به بدنه پروسسور یا همون پروسسور عضویت می‌دهد  
این عملی برای این که مشخص کند update یا reset انجام شود از دودار استفاده می‌کند



$$\hat{C}_t = \tanh(w_k[C_{t-1}, x_t] + b_c)$$

$$\pi_u = \sigma(w_2[c_{t-1}, x_t] + b_u)$$

$$C_t = \Gamma_u * \hat{C}_t + (1 - \Gamma_u) C_{t-1}$$

$$y_t = \text{softmax}(w_3 c_t + b_0)$$

اما نکته LSTM ~~اینست~~ یک نوع شبکه با قابلیت بلند مدت هست که مشکل vanishing gradient را حل کرد. در حقیقت این شبکه می تواند اطلاعات در رابطه خوبی یاد بگیرد.

۱. تفاوت اصلی GRU و LSTM نیز این است که GRU قابلیت update, reset را دارد ولی LSTM قابلیت ~~دری~~ و خروجی و فراموشی را دارد.

۲. ✓ در LSTM گیت  $\Gamma$  را داریم که شخص می شود چه مقدار از حافظه ~~در~~ state نگهدارنده شود و اندک ضعیف باشد به طریقی اندک حافظه از این خواهد رفت.

۳. ✓ اگر مجموع داده ها کوچک باشد GRU بهتر عمل خواهد کرد و LSTM با داده های زیاد بهتر کار می کند.

۴. ✓ GRU کل حافظه را در معیض نشان می دهد و LSTM این کار را انجام نمی دهد. و وضعیت های پنهان

استفاده از توابع activation مختلف دیدیم مانند Leaky-ReLU یا ReLU استفاده می کنیم (مانند randomized ReLU)  
این توابع به جای Sigmoid و tanh می توانند به کار روند ✓

✓ تدریس یا آموزش را گاهی هم : اگر تدریس خیلی بالا باشد ممکن است diverge کند و اگر تدریس خیلی پایین باشد با این که شبکه دیرتر آموزش پیدا می کند اما مشکل نه این (افزایش / کاهش بیش از حد) حل می شود.

✓ از مقداردهی اولیه وزن ها استفاده کنیم ( weight initialization ) : مثلا آموزش Xavier استفاده کنیم  
البته به خودی خود خوب نیست اولی با ترکیب روش های دیگر مدل robust تر می شود.

✓ از نرمال سازی بچ استفاده کنیم ( batch normalization ) : با این روش حتی اگر از tanh و Sigmoid استفاده کنیم مدل کمتر دچار مشکل می شود.

✓ انتخاب شبکه یا Architecture را تعیین کنیم : مثلا به جای معادله های قریبی تر از معادله های جدیدتر استفاده کنیم.