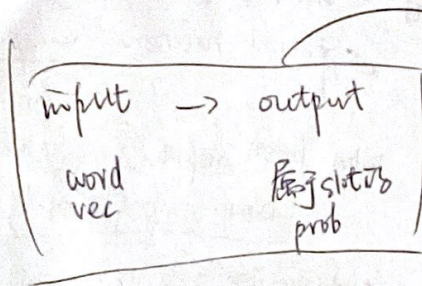


RNN

slot filling

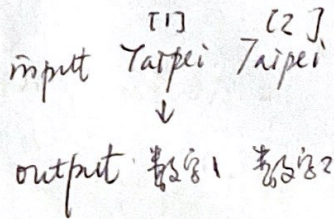
destination

slot { time of arrival

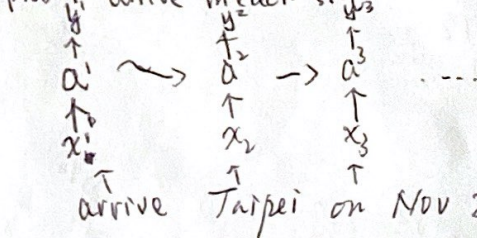


problem \rightarrow NN needs memory
 \downarrow
 RNN
 arrive Taipei on November 2nd
~~arrive~~ dest
 leave Taipei on November 2nd
 place of departure.

the output of hidden layer for memory

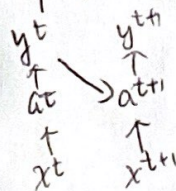
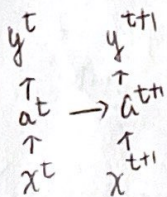


RNN. change the sequence order will change the output

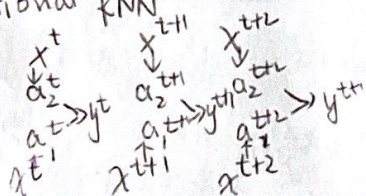


不是3个 network. 是同一个 network 在不同的时间被使用 3次.

Elman network & Jordan network (good performance)



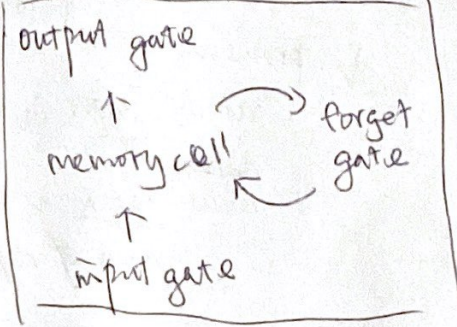
Bidirectional RNN



相对看到 x^t 之前和之后的才预测

LSTM

RNN 是真正的 short-term. 每一个 input 进来, memory 的值会被更新



special neuron: 4 inputs, 1 output
 z, z_i, z_f, z_o (z_i, z_f, z_o 是 $f(z)$)

why both input & output

\therefore input gate 可以控制是否存入 memory.

activation 函数是 sigmoid function $f \in (0,1)$
 for gates

z : input.

z_i : input gate

c : memory

z_f : forget gate

z_o : output gate

$z_i = Wz$ 0 代表遗忘, 1 代表记得.

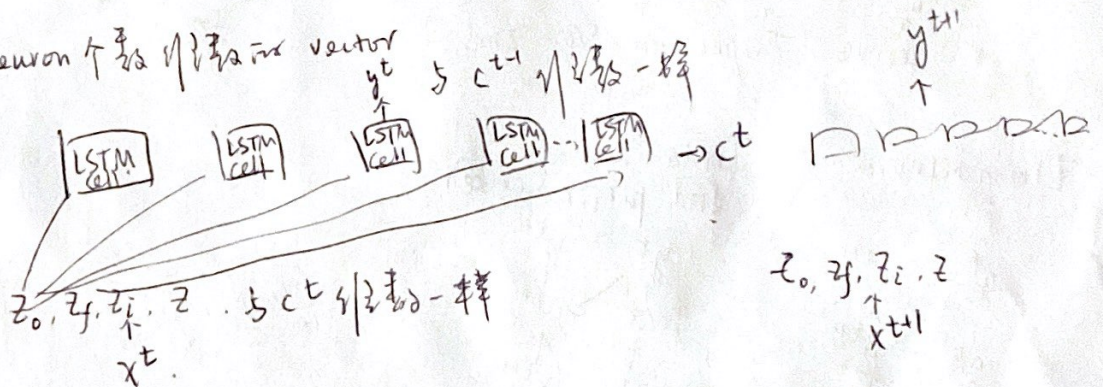
new memory: $c' = g(z) f(z_i) + c f(z_f)$

activation function

output: $a = h(c') f(z_o)$

c^{t-1} 是上层 neuron 个数的隐藏层 vector

$c^{t-1} : 1 \times 5$



实际上 x^t 是隐藏层 input word vector 以及 $h(c^{t-1})$. c^{t-1} (concatenate) peephole

keras supports: LSTM, GRU, Simple RNN layers.

3-4 gate 但 performance 差不多

update gate $z_t = f(W_z x_t + U_z h_{t-1} + b_z)$
 reset gate $r_t = f(W_r x_t + U_r h_{t-1} + b_r)$
 $\hat{h}_t = \phi(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h)$
 $h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \hat{h}_t$

BPTT (backpropagation through time). \therefore RNN 随着时间变化. $y_t = \sigma(h_t)$

RNN training

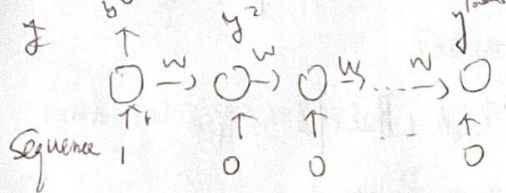
问题

1. The error surface is rough.

Solution - Clipping.

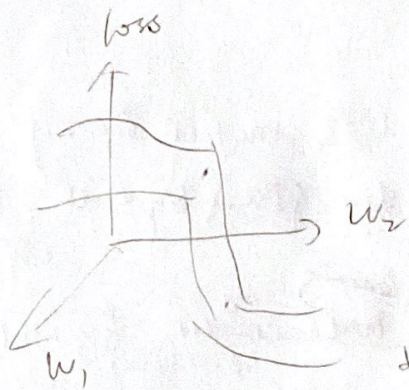
why

Toy example



$$w = 1.01, y^{1000} = 20000$$

$$w = 0.99, y^{1000} = 0$$



如果是在, w gradient 不为 0
para update 就不会很通

对于 sequence 长, 同一个 weight 在不同时间上反复被使用, explode or
LSTM 防止 gradient vanishing. 防止防止 gradient explode vanish

learning rate 比较小,

why? RNN memory 会遗忘 (→ 防止 forget)

LSTM memory 是保持的 防止 forget gate closed → No gradient vanishing if
forget gate is opened

GRU 只有 1 gate. 更加 robust. 如果 LSTM over fitting 很严重可以试 GRU.

GRU input gate 和 forget gate 联动起来. 当 input gate 打开, forget gate 会自动关闭 forget

Clockwise RNN is Structurally constrained Recurrent network (SCRN)

Vanilla RNN initiated with Identity matrix + ReLU activation performance is less

输出到输入 output 与 input 同时输入 RNN. no slot filling

但 RNN 可以做 many to one. many to many. In sentiment analysis document → embedding layer → output
attention feedforward network

RNN

* beyond sequence.

* syntactic parsing (结构 learning)

把解析树 parsing tree 转成 sequence (S(NP NNP)_{NP} (VP VPZ(NP DT IN))_{NP})_S
LSTM 可以处理, 不需要标注号.

* sequence to sequence (learning) auto-encoder 将 word sequence order 的情况作为 document embedding

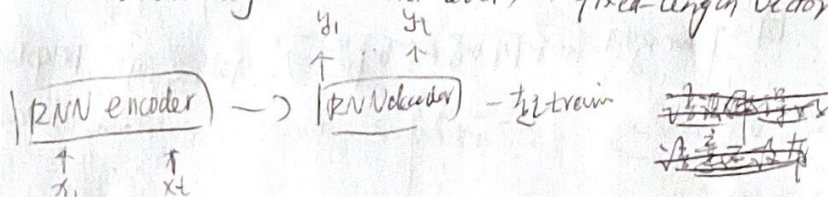
输入原文 \rightarrow encoder \rightarrow decoder \rightarrow 原文, 可以成功解码.

by document ~~embed~~ vector 代表了重要资讯 (且这种向量不需要 label data)

* skip-thought 是 \downarrow 类似 output 一个句子或短语, 能更好地理解语义.

* 也可以作 hierarchy. word vec \rightarrow sentence vec \rightarrow document vec \rightarrow sentence vec \rightarrow word vec

* dimension reduction for a sequence with variable length
audio segments (word-level) \rightarrow fixed-length vector



Attention-based model \rightarrow reading comprehension

visual QA

speech QA

RNN vs structured learning \rightarrow integrate deep. to constraint the iteration algo.

bi-directional LSTM + CRF / structured SVM for tagging

output
as feature