

No.

Date.

chap-6 循环神经网络

循环神经网络 RNN

递归神经网络

图网络

No.

Date.

6.1 神经网络增加记忆能力

短期记忆能力

6.1.1 延时神经网络 TDNN

时间维度上共享权重

延时单元

$$\vec{h}_t^{del} = f(\vec{h}_t^{del}, \dots, \vec{h}_{t-p+1}^{del})$$

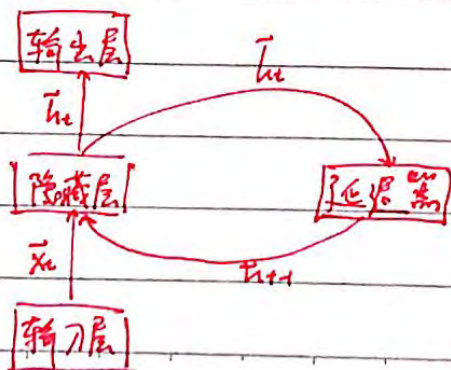
6.1.2 有外部输入的非线性自回归模型

自回归模型, AR

6.1.3 循环神经网络

$$\vec{x}_{1:T} = (\vec{x}_1, \vec{x}_2, \dots, \vec{x}_T) \quad \text{输入序列}$$

$$\vec{h}_t = f(\vec{h}_{t-1}, \vec{x}_t) \quad \text{隐状态}$$



6.2 简单循环网络

SRV

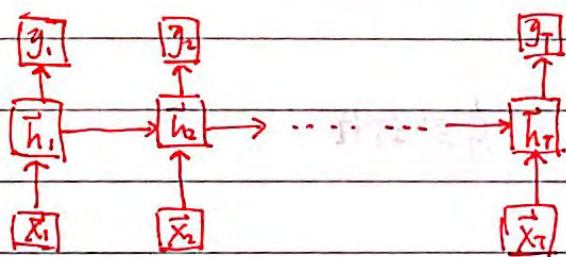
$$\bar{z}_t = U \bar{h}_{t-1} + W \bar{x}_t + \bar{b}$$

$$\bar{h}_t = f(\bar{z}_t)$$

U 状态到状态权重矩阵

W 状态-输入权重矩阵

\bar{z}_t 隐藏层的净输入



6.2.1 循环神经网络的计算能力

$$\bar{h}_t = f(U \bar{h}_{t-1} + W \bar{x}_t + \bar{b}) \quad (6.8)$$

$$\bar{y}_t = V \bar{h}_t \quad (6.9)$$

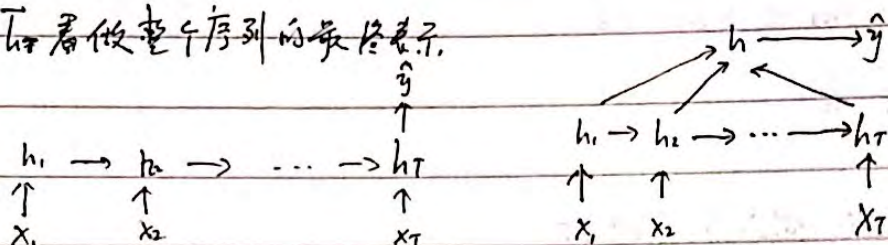
⇒ 通用近似定理

⇒ 图灵完备

No.

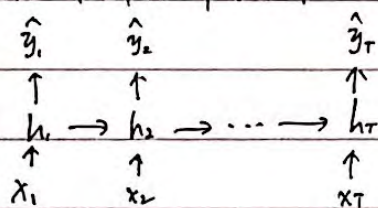
Date.

6.3 应用到机器学习

6.3.1 序列到类别模型 分类将 \vec{h}_T 看做整个序列的最终表示。

$$\hat{y} = g(\vec{h}_T)$$

6.3.2 同步的序列到序列模型



序列标注

6.3.3 异步的序列到序列模型

编码器-解码器模型。

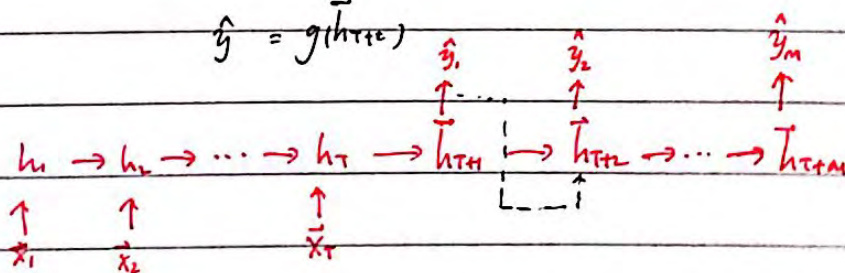
$$\vec{x}_{1:T} = (\vec{x}_1, \dots, \vec{x}_T)$$

$$\vec{y}_{1:m} = (\vec{y}_1, \dots, \vec{y}_m)$$

$$\vec{h}_t = f_1(\vec{h}_{t-1}, \vec{x}_t)$$

$$\vec{h}_{t+m} = f_2(\vec{h}_{t+m-1}, \vec{y}_t)$$

$$\hat{y} = g(\vec{h}_{t+m})$$



6.4 参数学习

$$L_t = L(y_t, \vec{g}_t)$$

t时刻的损失函数

$$g_t(\vec{h}_t)$$

t时刻的输出

$$L = \sum_t L_t$$

整个序列上的损失函数

$$\frac{\partial L}{\partial U} = \sum_t \frac{\partial L_t}{\partial U} \leftarrow \vec{h}_t = f(U\vec{h}_{t-1} + W\vec{x}_t + b)$$

随时间反向传播 BPTT

实时循环学习算法 RTRL

6.4.1 随时间反向传播算法

BPTT将RNN看成一个展开的多层前馈网络。每一层对应每个时刻。

在展开的网络中所有层的参数是共享的，因此参数的真实梯度是将所有展开层的参数梯度之和。

⇒ 计算偏导数 $\frac{\partial L_t}{\partial U}$

$$\vec{s}_k = U\vec{h}_{k-1} + W\vec{x}_k + b \quad \text{时刻k时的净输入。}$$

$$\frac{\partial L_t}{\partial U_{ij}} = \sum_k \left(\left(\frac{\partial L_t}{\partial \vec{s}_k} \right)^T \frac{\partial^+ \vec{s}_k}{\partial U_{ij}} \right)$$

$$= \sum_k \left(\left(\frac{\partial^+ \vec{s}_k}{\partial U_{ij}} \right)^T \frac{\partial L_t}{\partial \vec{s}_k} \right)$$

why

(6.3.4)

No.

Date.

$\frac{\partial \bar{z}_h}{\partial u_{ij}}$ 表示直接导数

$$\frac{\partial \bar{z}_h}{\partial u_{ij}} = \begin{bmatrix} 0 \\ [\bar{h}_{h+1}]_j \\ 0 \end{bmatrix} \doteq \mathbb{I}_h([\bar{h}_{h+1}]_j) \quad (6.30)$$

$\delta_{t,h} = \frac{\partial L_t}{\partial \bar{z}_h}$ 定义为 t 时刻的损失对 h 时刻隐藏神经层的净输入 \bar{z}_h 的导数

$$\begin{aligned} \delta_{t,h} &= \frac{\partial L_t}{\partial \bar{z}_h} \\ &= \frac{\partial \bar{h}_h}{\partial \bar{z}_h} \cdot \frac{\partial \bar{z}_{h+1}}{\partial \bar{h}_h} \cdot \frac{\partial L_t}{\partial \bar{z}_{h+1}} \\ &= \text{diag}(f'(\bar{z}_h)) \cdot U^T \cdot \delta_{t,h+1} \end{aligned} \quad (6.31)$$

$$\bar{z}_t = U\bar{h}_{t+1} + W\bar{x}_t + \bar{b}$$

$$\bar{h}_t = f(\bar{z}_t) \quad (6.32)$$

将式 (6.32) $\delta_{t,h}$ 和 (6.31) $\frac{\partial \bar{z}_h}{\partial u_{ij}}$ 代入 (6.30) 得到 $\frac{\partial L_t}{\partial u_{ij}}$

$$\frac{\partial L_t}{\partial u_{ij}} = \sum_k [\delta_{t,h}]_k [\bar{h}_{h+1}]_j$$

$$= \bar{z} \otimes \delta_{t,h} \bar{h}_{h+1}^T$$

⇒ 多数梯度

$$\frac{\partial L}{\partial U} = \sum_{t=1}^T \sum_{h=1}^H \delta_{t,h} \vec{h}_{t-1}^T$$

$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \sum_{h=1}^H \delta_{t,h} \vec{x}_{t-1}^T$$

$$\frac{\partial L}{\partial b} = \sum_{t=1}^T \sum_{h=1}^H \delta_{t,h}$$

⇒ 计算复杂度

一个完整的前向计算和反向计算。

6.4.2 实时循环学习算法

$$\frac{\partial L_t}{\partial U_{ij}} = \left(\frac{\partial h_{t+1}}{\partial U_{ij}} \right)^T \cdot \frac{\partial L_t}{\partial h_t}$$

⇒ 两种算法比较

BPTT: 一般网络输出维度 < 输入维度, BPTT算法更适用. 空间复杂度低.

TRL: 有限序列, 无限序列.

No.

Date.

6.5 长期依赖问题

在 BPTT 算法中, 将公式 (6.34) 展开, 得到

$$\delta_{t,h} = \frac{1}{\pi} \sum_{i=h}^{t-1} (\text{diag}(f'(\vec{z}_i)) U^T) \cdot \delta_{t,i}$$

$$\frac{1}{2} r \approx \|\text{diag}(f'(\vec{z}_i)) U^T\|$$

$$\delta_{t,h} = \frac{1}{\pi} r^{t-h} \delta_{t,t}$$

$r > 1$, 梯度爆炸

$r < 1$, 梯度消失. 无法影响远处信息.

长期依赖问题

6.5.1 改进方案

梯度爆炸

权重衰减, 对范数正则化, 使 $r < 1$

梯度截断

梯度消失

改变模型

让 $U = I$, $f'(\vec{z}_i) = 1$, 即

$\vec{h}_t = \vec{h}_{t-1} + g(\vec{x}_t; \theta)$ 丢失了非线性性质

增加记忆容量

$$\vec{h}_t = \vec{h}_{t-1} + g(\vec{x}_t, \vec{h}_{t-1}; \theta)$$

6.6 基于门控的循环神经网络.

1. 增加一些额外的存储单元, 外部记忆
2. 进行选择性遗忘, 选择性更新

6.6.1 ~~长期~~ 长短期记忆网络 (LSTM)

新的内部状态 c_t (记忆单元)

c_t 专门进行线性的循环信息传递, 同时非线性输出信息给隐藏层的外部状态 h_t

$$\vec{c}_t = \vec{f}_t \odot \vec{c}_{t-1} + \vec{i}_t \odot \vec{c}_t \quad \text{线性}$$

$$\vec{h}_t = \vec{o}_t \odot \tanh(c_t) \quad \text{非线性}$$

\vec{c}_t 假状态.

$$\vec{c}_t = \tanh(w_c \vec{x}_t + v_c \vec{h}_{t-1} + \vec{b}_c)$$

门机制

取值在 (0, 1) 之间.

遗忘门 \vec{f}_t

\vec{c}_t 需要遗忘多少信息

输入门 \vec{i}_t

\vec{c}_t 有多少信息要保存.

输出门 \vec{o}_t

\vec{c}_t 有多少信息传给 \vec{h}_t

$\vec{f}_t = 0, \vec{i}_t = 1$ 时 历史信息被遗忘

$\vec{f}_t = 1, \vec{i}_t = 0$ 时. 复制历史信息

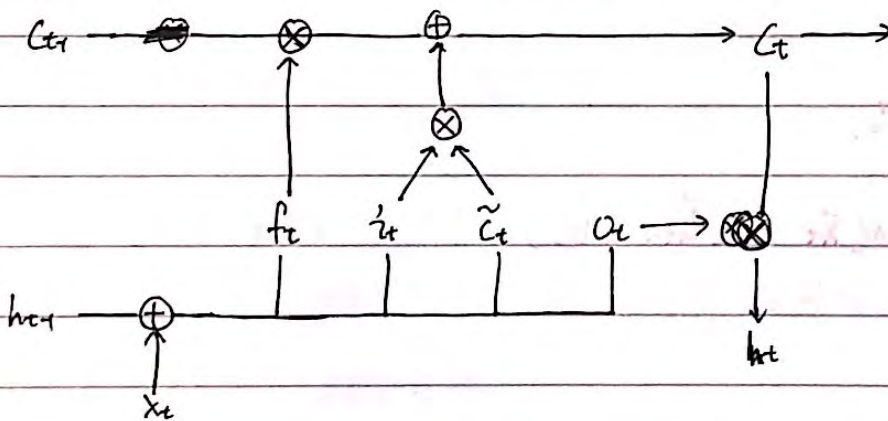
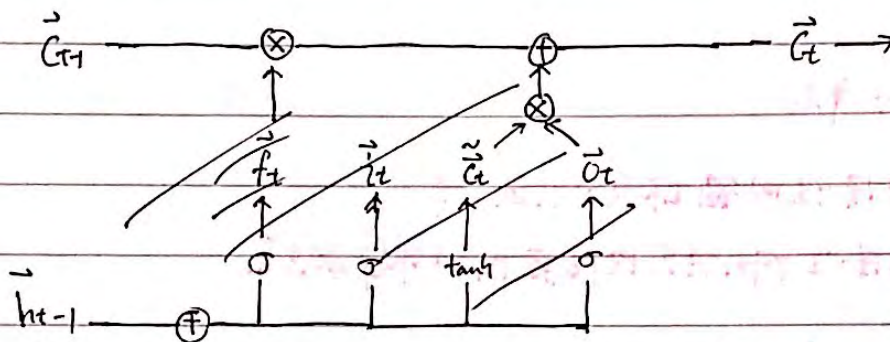
No.

Date.

① $\vec{h}_{t-1} \vec{x}_t \Rightarrow \vec{f}_t, \vec{i}_t, \vec{\tilde{c}}_t, \vec{o}_t$ (PI State)

② $\vec{f}_t \vec{x}_t, \vec{i}_t, \vec{\tilde{c}}_t \Rightarrow \vec{c}_t$

③ $\vec{o}_t \Rightarrow \vec{h}_t$



$$\begin{bmatrix} \vec{\tilde{c}}_t \\ \vec{o}_t \\ \vec{i}_t \\ \vec{f}_t \end{bmatrix} = \begin{bmatrix} \tanh \\ \sigma \\ \sigma \\ \sigma \end{bmatrix} (W \begin{bmatrix} \vec{x}_t \\ \vec{h}_{t-1} \end{bmatrix} + b)$$

$$\vec{c}_t = \vec{f}_t \otimes \vec{c}_{t-1} + \vec{i}_t \otimes \vec{\tilde{c}}_t$$

$$\vec{h}_t = \vec{o}_t \otimes \tanh(\vec{c}_t)$$

⇒ 记忆

短期记忆 h_t

长期记忆 c_t

6.6.2 LSTM 网络的各种变化

无遗忘门的 LSTM 网络

peephole 连接

耦合输入门和遗忘门

6.6.3 门控循环单元网络 (GRU)

输入门
遗忘门

更新门

\bar{h}_t 和 \bar{h}_{t-1} 之间引入线性依赖关系

$h_{t-1} \rightarrow \otimes$

$\rightarrow \oplus$

No.

Date.

救世主是由已知和设计师的状态决定的。

救世主

$$\tilde{h}_c = \tanh(W_c \tilde{x}_c + U(r_c \odot h_{c-1}) + b_h)$$

↑
救世主

↑
已知

↑
设计师

↑
历史

设计师控制上一代救世主传递多少信息给当前代

设计师

$$r_c = \sigma(W_r \tilde{x}_c + U_r h_{c-1} + b_r)$$

↑
已知

↑
上一代救世主

救世主的信息

$$h_c = z_c \odot h_{c-1} + (1 - z_c) \odot \tilde{h}_c$$

↑
已知

↑
上一代救世主

↑
更新数据

↑
救世主

已知影响上一代救世主

更新数据影响救世主

已知和设计师同时影响救世主

已知

$$z_c = \sigma(W_z \tilde{x}_c + U_z h_{c-1} + b_z)$$

↑
已知

↑
上一代救世主

$$r_c = \sigma(W_r \tilde{x}_c + U_r h_{c-1} + b_r)$$

$$z_c = \sigma(W_z \tilde{x}_c + U_z h_{c-1} + b_z)$$

$$\tilde{h}_c = \tanh(W_c \tilde{x}_c + U(r_c \odot h_{c-1}) + b_h)$$

$$h_c = z_c \odot h_{c-1} + (1 - z_c) \odot \tilde{h}_c$$

重置门

$$r_c = 0 \quad \text{在历史数据}$$

$$= 1 \quad \text{和历史数据交互 SRN}$$

更新门

$$z_c = 0, r_c = 1 \quad \text{SRN}$$

6.7 深层循环神经网络.

增加隐藏状态到输出 $\vec{h}_t \rightarrow \vec{y}_t$

输入到隐藏状态 $\vec{x}_t \rightarrow \vec{h}_t$

6.7.1 堆叠循环神经网络 (SRNN)

$$\vec{h}_t^{(d)} = f(U^{(d)} \vec{h}_{t-1}^{(d)} + W^{(d)} \vec{x}_t^{(d)} + b^{(d)})$$

6.7.2 双向循环神经网络 (Bi-RNN)

$$\vec{h}_t^{(u)} = f(U^{(u)} \vec{h}_{t-1}^{(u)} + W^{(u)} \vec{x}_t + b^{(u)})$$

$$\vec{h}_t^{(v)} = f(U^{(v)} \vec{h}_{t-1}^{(v)} + W^{(v)} \vec{x}_t + b^{(v)})$$

$$\vec{h}_t = \vec{h}_t^{(u)} \oplus \vec{h}_t^{(v)}$$

No.

Date.

6.8 扩展到图结构.

每个节点都收到其邻居的输入.

6.8.1 递归神经网络 (RecNN)