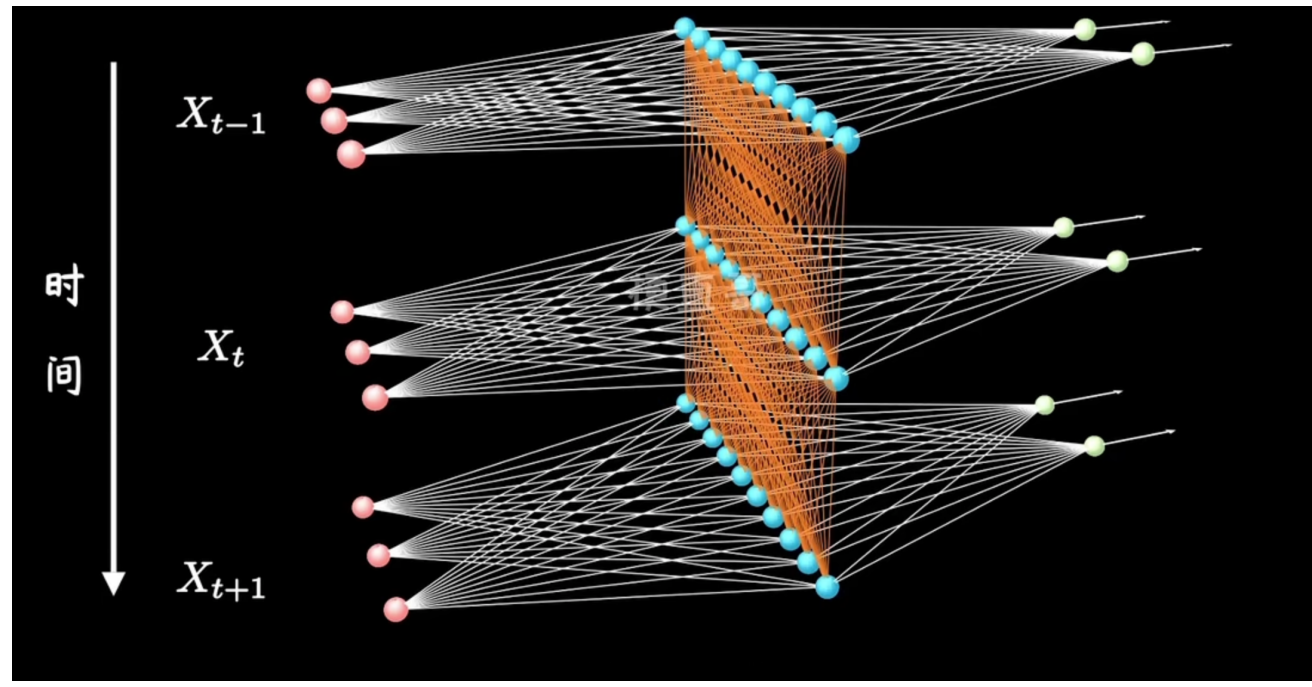


循环神经网络

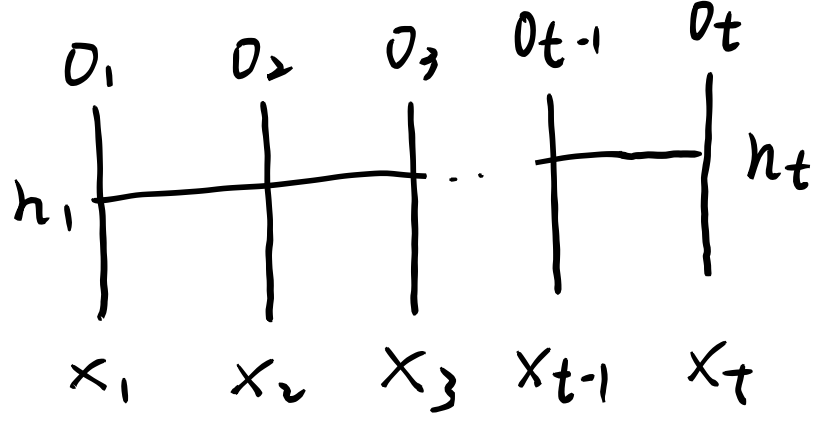


在样本之间具有关联性

例如 图中的三层神经网络 $h: size(1 \times 5)$ $o: size(1 \times 1)$
 令 单例 $x: size(1 \times n)$ $w_{xh}: size(n \times 5)$ $w_{hh}: size(5 \times 5)$
 $h_t = w_{xh} x_t + w_{hh} h_{t-1} + b_h$
 $o = w_{ho} \sigma(h_t) + b_o$

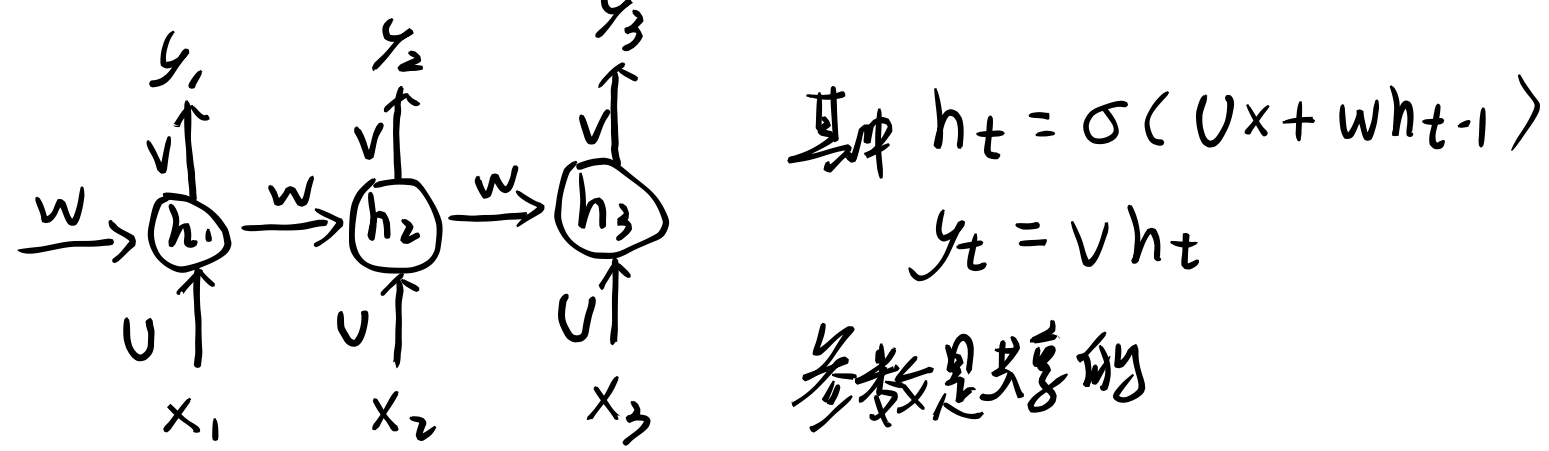
反向传播:

定义损失函数:



$$J = J_1(o_1) + J_2(o_2) + \dots + J_t(o_t)$$

一个例子: 看其中细偏置元



$$h_t = \sigma(Ux_t + wh_{t-1})$$

$$y_t = Vh_t$$

参数是共享的

$$h_1 = \sigma(Ux_1 + wh_0) \quad h_2 = \sigma(Ux_2 + wh_1) \quad h_3 = \sigma(Ux_3 + wh_2)$$

$$y_1 = Vh_1 \quad y_2 = Vh_2 \quad y_3 = Vh_3$$

$$L = \sum L_i$$

在 t=3 时求对 U, w, v 求偏导

$$\frac{\partial L_3}{\partial v} = \frac{\partial L_3}{\partial y_3} \frac{\partial y_3}{\partial v}$$

$$\frac{\partial L_3}{\partial U} = \frac{\partial L_3}{\partial y_3} \frac{\partial y_3}{\partial h_3} \left(\frac{\partial h_3}{\partial U} + \frac{\partial h_3}{\partial h_2} \left(\frac{\partial h_2}{\partial U} + \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial U} \right) \right) = \sum_{k=1}^t \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial U} = \sum_{k=1}^t \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial U}$$

$$\frac{\partial L_3}{\partial w} = \frac{\partial L_3}{\partial y_3} \frac{\partial y_3}{\partial h_3} \left(\frac{\partial h_3}{\partial w} + \frac{\partial h_3}{\partial h_2} \left(\frac{\partial h_2}{\partial w} + \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w} \right) \right) = \sum_{k=1}^t \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial w}$$

$$\Rightarrow \frac{\partial L}{\partial v} = \sum_{t=1}^T \frac{\partial L_t}{\partial v} = \sum_{t=1}^T \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial v}$$

$$\frac{\partial L}{\partial U} = \sum_{t=1}^T \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \sum_{k=1}^t \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial U}$$

$$\frac{\partial L}{\partial w} = \sum_{t=1}^T \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \sum_{k=1}^t \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial w}$$

代入

$$\frac{\partial L}{\partial U} = \sum_{t=1}^T \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \sum_{k=1}^t (\sigma'w)^{t-k} \frac{\partial h_k}{\partial U}$$

若 $\sigma'w < 1$ 当 t-k 过大时 $\sigma'w \rightarrow 0$ 导致梯度消失

若 $\sigma'w > 1$ 当 t-k 过大时 $\sigma'w \rightarrow \infty$ 导致梯度爆炸

$$(Rule) \sigma' = \text{diag}(0, 1, \dots)$$

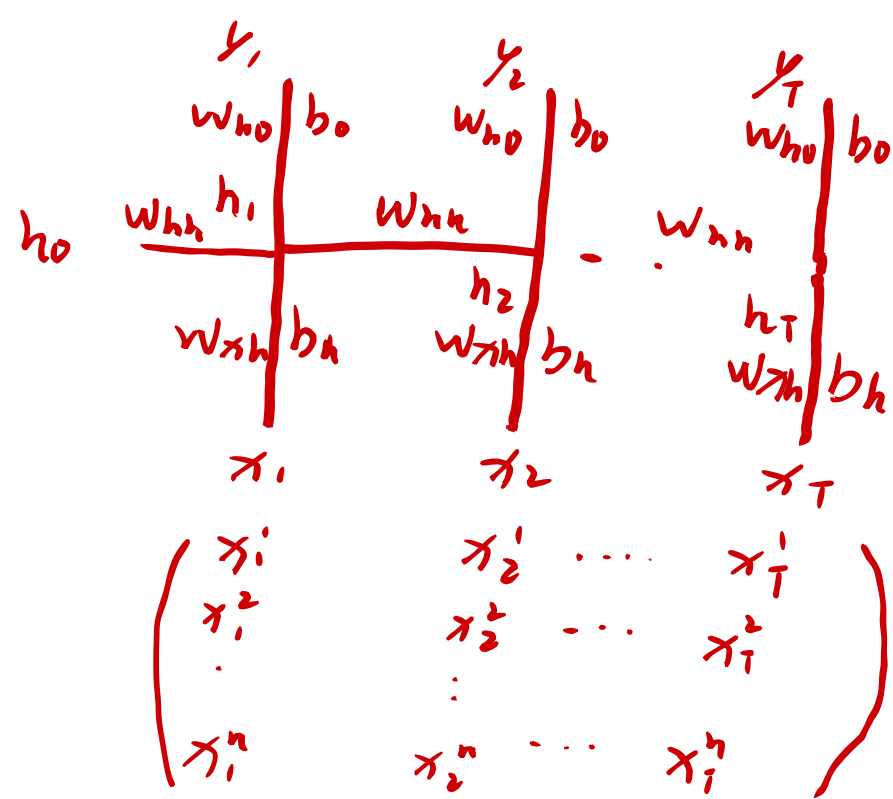
其中 $\frac{\partial h_j}{\partial h_{j-1}} = \sigma'w$ 故 $\frac{\partial h_t}{\partial h_k} = (\sigma'w)^{t-k}$
 相当于在 w 中某一行置为 0 (初等行变换)

$$y = \sigma(x) \text{ 其中 } x, y \in \mathbb{R}^n \text{ 则 } \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \sigma \left(\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \right)$$

$$\frac{\partial y}{\partial x} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \frac{\partial y_n}{\partial x_2} & \dots & \frac{\partial y_n}{\partial x_n} \end{pmatrix} = \text{diag}(0, 1, 0, \dots)$$

$$\begin{cases} \frac{\partial g(u(x))}{\partial x} = \frac{\partial g}{\partial u} \frac{\partial u}{\partial x} \text{ 其中 } g, u, x \text{ 都是向量.} \\ (Ax)' = A \quad A \text{ 是矩阵 } x \text{ 是向量.} \end{cases}$$

从直观现代的解释



(batch-size x T) 张量的网络输入

定义的 RNN 网络只有三层即 输入层, 隐变量层, 输出层

上图所示是 RNN 的一个时间步的网络

上图所示在神经网络中对应的是一个字母的 one-hot 编码

因此 RNN input-size 为 28, output-size 也是 28 从而输出也对应一个字母

开始训练时 采用批量训练

for x, y in train_iter:
 y_hat = net(x)

因此代码中 h0 为一个大小为 (batch-size, hidden-size) 的合理矩阵

这本身 (h_1, h_2, \dots, h_T) 还是一个尺寸为 (时间步, batch-size, hidden-size) 的张量

最后输出的 (y_1, y_2, \dots, y_T) 的形状为 (时间步, batch-size, 28) (28 是词表大小)

而这里的 y 的形状是 (batch-size, 时间步)

因此计算 loss 时需要转置

即:

$$y = y.T.reshape(-1)$$

$$y_hat = y_hat.reshape(-1, 28)$$

$$loss = loss(y_hat, y)$$