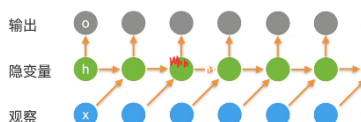


# RNN

- 潜变量自回归模型：使用潜变量 $h_t$ 总结过去信息
- 循环神经网络：



$$p(h_t | h_{t-1}, x_{t-1}) ; p(x_t | h_t, x_{t-1})$$

$$h_t = \varphi(W_{hh} h_{t-1} + W_{hx} x_{t-1} + b_h)$$

→ 去掉就退化为MLP

$$O_t = W_{ho} h_t + b_o$$

- 困惑度 (Perplexity)

- 衡量一个语言模型好坏可以用平均交叉熵
- 历史原因NLP使用困惑度 $\exp(\pi)$ 来衡量，1表示完美， $\infty$ 是最差情况

$$\pi = \frac{1}{n} \sum_{i=1}^n -\log p(x_t | x_{t-1} \dots)$$

n为seq长度

- 梯度剪裁

- 迭代中计算T个时间步上的梯度，反向传播产生 $O(T)$ 的矩阵乘法，导致数值不稳定
- 梯度剪裁能有效预防梯度爆炸，如果提督长度超过 $\theta$ ，那么拖回长度 $\theta$
- 这样做后梯度范数不会超过 $\theta$ ，且更新后的梯度仍与g原始方向对齐

$$g' \leftarrow \min(1, \frac{\theta}{\|g\|}) \cdot g$$

g为所有层梯度

$\|g\|$ 代表梯度长度( $L_2$ 范数), 超过 $\theta$ 就压缩

相当于将每个分量压缩了

也为模型赋予了稳定性

