a language model that can handle unbounded context (prefix).

$$p(X_1, X_2, \cdots, X_T) = \prod_{t=1}^{T} p(X_t | X_1, \cdots, X_{t-1})$$

$X_t \in V$

$$\approx p(X_t | X_{t-n}, \cdots, X_{t-1}) : n\text{-gram LM}$$

---

Neural probabilistic language model.

① $\underbrace{e(X_{t-n}, \cdots, X_{t-1})}_{} = [e(x_{t-n})^T ; \cdots , e(x_{t-1})^T]^T \in \mathbb{R}^{nd}$

② $h = \tanh(W \underbrace{e(x_{t-n}, \cdots, x_{t-1})}_{} + b) \in \mathbb{R}^{d'}$

$W \in \mathbb{R}^{d' \times nd}$

③ $F(x_t, x_{t-n:t-1}) = e'(x_t)^T \underset{\uparrow}{h} + b'(x_t)$

$W \in \mathbb{R}^{d' \times nd}$    classembedding

$$
\begin{array}{c}
d' 
\begin{bmatrix} W_{t-n} \end{bmatrix}
\begin{bmatrix} W_{t-n+1} \end{bmatrix}
\cdots
\begin{bmatrix} W_{t-1} \end{bmatrix}
\begin{bmatrix} e(X_{t-n}) \\ e(X_{t-n+1}) \\ \vdots \\ e(X_{t-1}) \end{bmatrix}
= \sum_{n=1}^{n} W_{t-n'} \, e(x_{t-n'})
\end{array}
$$

$\underbrace{\quad}_{d} \quad \underbrace{\quad}_{d} \quad \underbrace{\quad}_{d}$

---

A text classifier summarizes the input text

$$X = (x_1, \cdots, x_T)$$

$$\phi(x) = \sigma\left(W \sum_{t=1}^{T} e(x_t) + b\right) \quad , \quad \boxed{F(x,y) = n_y^T \phi(x) + c_y}$$
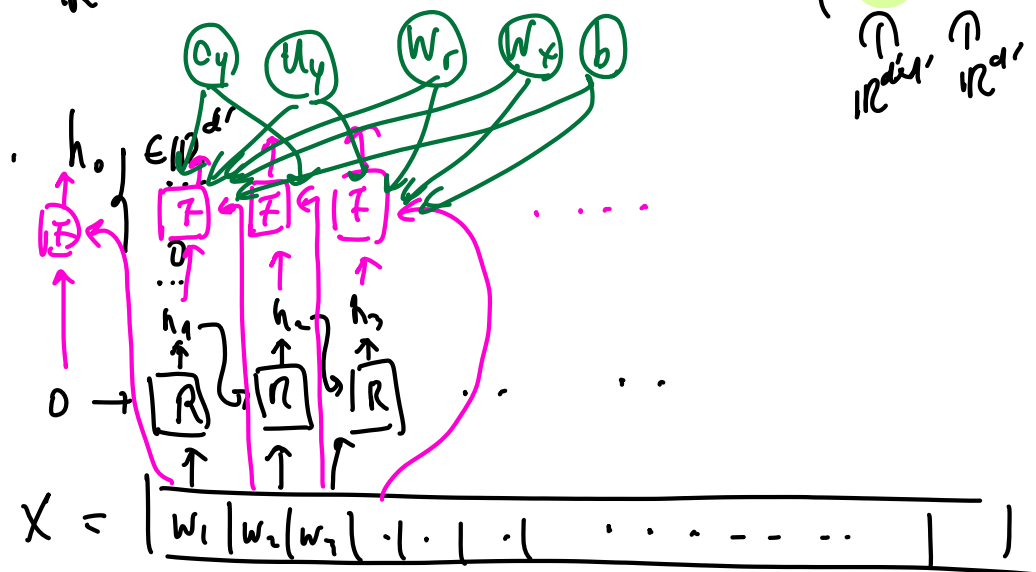
A recurrent neural network.

  ∴ recursively compress the input sequence.

If $h_{t-1} = \phi(x = [w_1, \cdots, w_{t-1}])$, can we add $w_t$?

· $h_t = R(h_{t-1}, w_t)$, where e.g. $R(h_{t-1}, w_t)$

$$= \sigma\left( W_r \, h_{t-1} + W_x \, e(w_t) + b \right)$$

$$\underset{\mathbb{R}^{d' \times d'}}{\Uparrow} \quad \underset{\mathbb{R}^{d'}}{\Uparrow} \quad \underset{\mathbb{R}^{d' \times d}}{\Uparrow} \quad \underset{\mathbb{R}^{d}}{\Uparrow} \quad \underset{\mathbb{R}^{d'}}{\Uparrow}$$
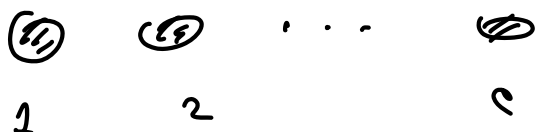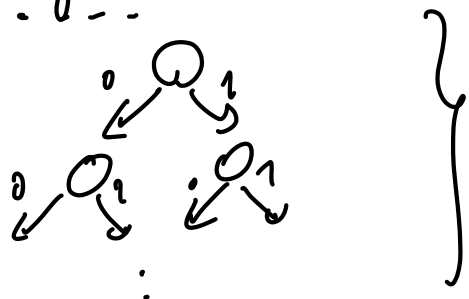


A recurrent neural net language model [Mikolov et al. 2010]

Alice $\implies$ Bob

$m \in \{1, \cdots, C\}$

$\log_2 C$ bits



$- \sum p(m) \log_2 p(m)$

$$2^{-\log_2 p(m)} \quad \Rightarrow \quad \boxed{2^{-\log_2 p(w_t \mid w_{<t})}}$$

$$B = -\frac{1}{\sum_{n=1}^{N} T^n} \sum_{n=1}^{N} \sum_{t=1}^{T^n} \log_2 p(w_t^n \mid w_{<t}^n) \quad , \quad \text{perplexity} = 2^B$$