We can adapt the weights between $V_0$ and $V_1$ so that for every $i \in [k]$, the $i^{th}$ neuron is $g_i(x)$ i

The neuron in the output layer implements a disjunction i.e. OR of the functions $g_i(x)$.

AND
Bias $-k+1$

$$f(x) = \text{sign}\left( \sum_{i=1}^{k} g_i(x) + k - 1 \right)$$

bias.

$-k + k - 1 \geq -1$

$-(k-2) + k + 1 \geq 1$

$\text{sign}(1)$

Even if we try to model functions of the form $\{0,1\}^n \to \{0,1\}$ the size of the network will be exponential in $n$.

A neural network can approximate 1-Lipschitz function $f: [-1, +1]^n \longrightarrow [-1, 1]$ within a precision $\epsilon$, but the size of the network will be exponential in $n$.

$f(u)$
$f(v)$

$n$-variables $[-1 \quad +1]$

$\|f(u) - f(v)\| \leq \|u - v\|^2$

$\to [-1, +1]$

Softmax converts $k$ real valued predictions $v_1 \cdots v_k$ into output probabilities $o_1 \cdots o_k$ using the relation

$v$
logits

$o$
softmax layer

probability

$$o_i = \frac{\exp(v_i)}{\sum_{j=1}^{k} \exp(v_j)} \qquad \forall i \in \{1, \ldots k\}$$

Not an element-wise operation

The softmax is mostly paired with the cross-entropy loss. If the target probability distribution over the k-classes is given by the vector $y_1 \cdots y_k$ GT then the cross-entropy loss is defined as

$$L = -\sum_{i=1}^{k} y_i \log(o_i)$$

loss.

$\underset{\text{output}}{o} \quad \underset{\text{GT}}{y}$

Ground truth $\underline{O}$

Labels 1 $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\underline{y}$ $\longrightarrow$ $\begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \end{bmatrix}$ $\begin{matrix} 0.05 \\ 0.7 \\ 0.1 \\ 0.1 \\ 0.05 \end{matrix}$ $\Big\} \sum P_i = 1.$

2

3

4

5

$$\boxed{Loss = -\log P_{GT}}$$
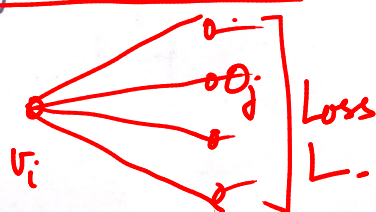
$\text{if } P_2 = 1 \quad P_2 \to 0$

$Loss = 0.$

$\underline{Loss = \infty, \text{ if } P_2 = 0}$

The partial derivative of the loss w.r.t. the logit value $v_i$

**Proof** is

$$\frac{\partial L}{\partial v_i} = \sum_{j=1}^{k} \left(\frac{\partial L}{\partial o_j}\right)\left(\frac{\partial o_j}{\partial v_i}\right) = o_i - y_i$$


Loss $L$.

$o_i = \dfrac{\exp(v_i)}{\sum_j \exp(v_j)}$   softmax

$$\frac{\partial o_j}{\partial v_i} = -\frac{\exp(v_i)\,\exp(v_j)}{\left(\sum \exp(v_i)\right)^2} + \frac{\exp(v_i)}{\sum_j \exp(v_j)}$$

$-o_i o_j$

$o_i$

the 2nd term is there if $i=j$

$$\frac{\partial o_j}{\partial v_i} = -o_i o_j + o_i \qquad \text{if } i=j$$
$$= o_i(1-o_i) \qquad \text{if } i=j$$

$$\frac{\partial o_j}{\partial v_i} = -o_i o_j \qquad \text{if } i \neq j$$

Substituting these partial derivatives

$$-\frac{\partial L}{\partial v_i} = \sum_{j=1}^{k} \left(\frac{y_i}{o_j}\right) o_j(1-o_j) \qquad \text{for } i=j \quad \frac{\partial o_j}{\partial v_i} = o_i(1-o_i)$$

$\frac{\partial o_j}{\partial v_i}$

$i \neq j \quad \frac{\partial o_j}{\partial v_i} = -o_i o_j \quad \text{if } i \neq j$

$$= \sum_{\substack{j=1 \\ j \neq i}}^{k} \frac{y_i}{o_j}(-o_i o_j) + \frac{y_j}{o_j} o_j(1-o_j)$$

$j=i$

$$= \sum_{\substack{j=1 \\ j \neq i}}^{k} -y_j o_i + y_j - y_j o_j$$
$j=i$   $j=i$

$= o_i - y_i$

$$= y_j - o_j \left(\sum_{j=1}^{k} y_j\right) = y_j - o_j$$
$j=i$   $j=i$   $=1$   $j=i$

$$\therefore \frac{\partial L}{\partial v_i} = o_j - y_j$$
$j=i$

# Recurrent Neural Network

$G_{V,E,\sigma}$

speech, time-series
sentences } text
handwriting }

The feedforward dense connected neural network is unaware
of the dependencies/orderings of the feature components of
the input. Though the network can discover such orderings,
it will require a very large training set.
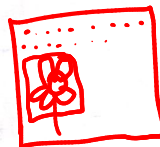For applications involving sequential input such as:
Language Modeling → predict the next word in a sentence
The sun rises...
Sentence Translation, sentence classification
sentiment
Image Captioning
— the input sequence length is not constant.

To exploit such local dependencies in sentences, images, speech,
time series data, the network architecture should be aware
of the ordering of the input.
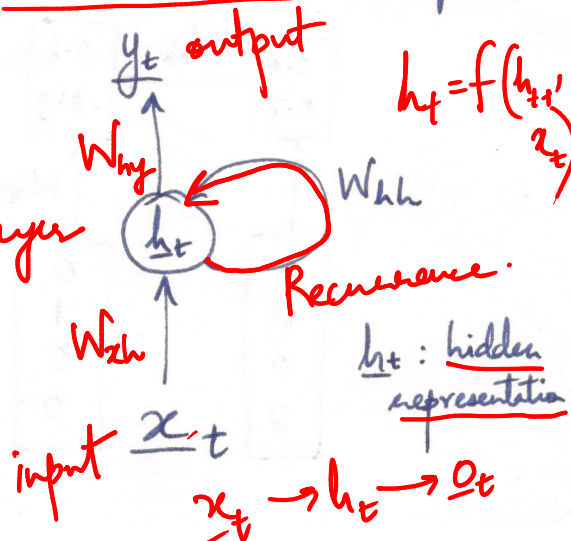repeated

① aware of ordering
② variable size of input

A recurrent neural network (RNN) uses a variable number of
layers such that each layer takes input of a specific position
(time stamp)
in a sequence. The inputs are provided in a sequence to the
"temporal layers". Each layer takes a multidimensional input
and produces a multi-dimensional output.

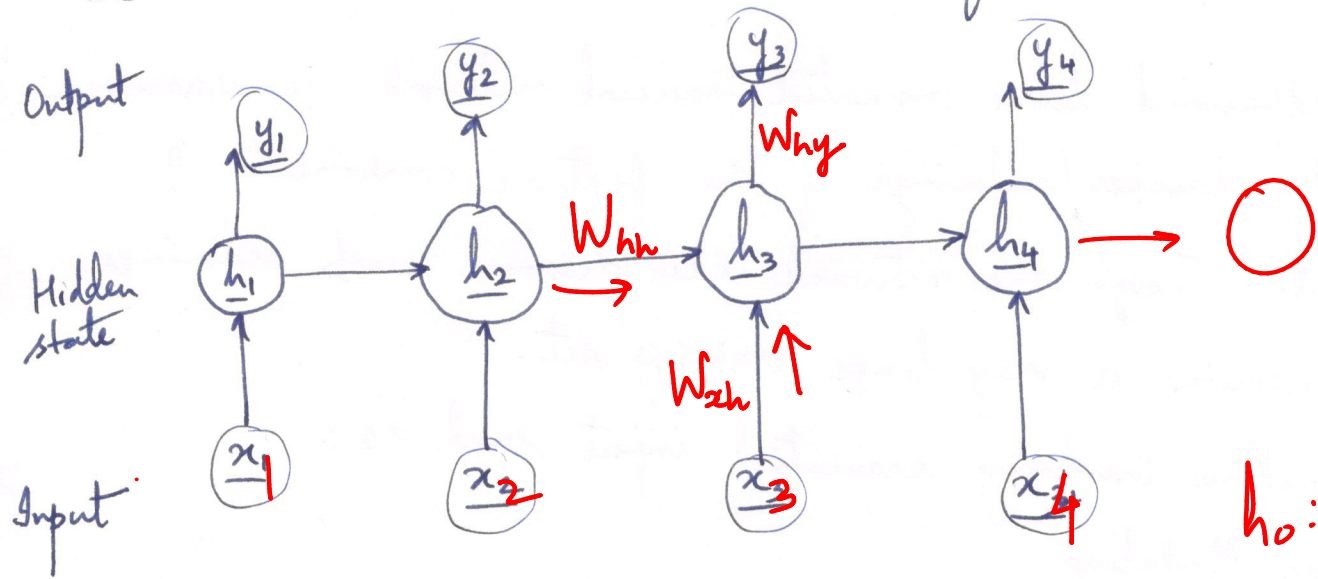Further, each layer uses the same
set of parameters (repeated single layer
recurrent computation)
This is called as parameter sharing.
Every input is treated in the same
manner.

$y_t$ output

$W_{hy}$     $W_{hh}$

$h_t$     Recurrence

$W_{xh}$

input $x_t$

$h_t = f(h_{t-1}, x_t)$

$h_t$ : hidden representation

$x_t \rightarrow h_t \rightarrow o_t$

The recurrent network can be unrolled in time

Output



$$h_t = f(h_{t-1}, x_t)$$

$$y_t = g(h_t)$$

$$h_1 = f(h_0, x_1)$$

$$h_2 = f(\underbrace{f(h_0, x_1)}_{h_1}, x_2)$$

$\vdots$

Since $y_t = g(h_t)$

$$\boxed{h_t = \underbrace{\tanh}_{\text{nonlinearity}} (\underbrace{W_{xh} x_t + W_{hh} h_{t-1}}_{\text{Affine}})}$$

$$y_t = W_{hy} h_t$$

$$\boxed{y_t = F_t(x_1, x_2, \ldots x_t)}$$