# Back-Propagation (The Generalized Delta Rule)
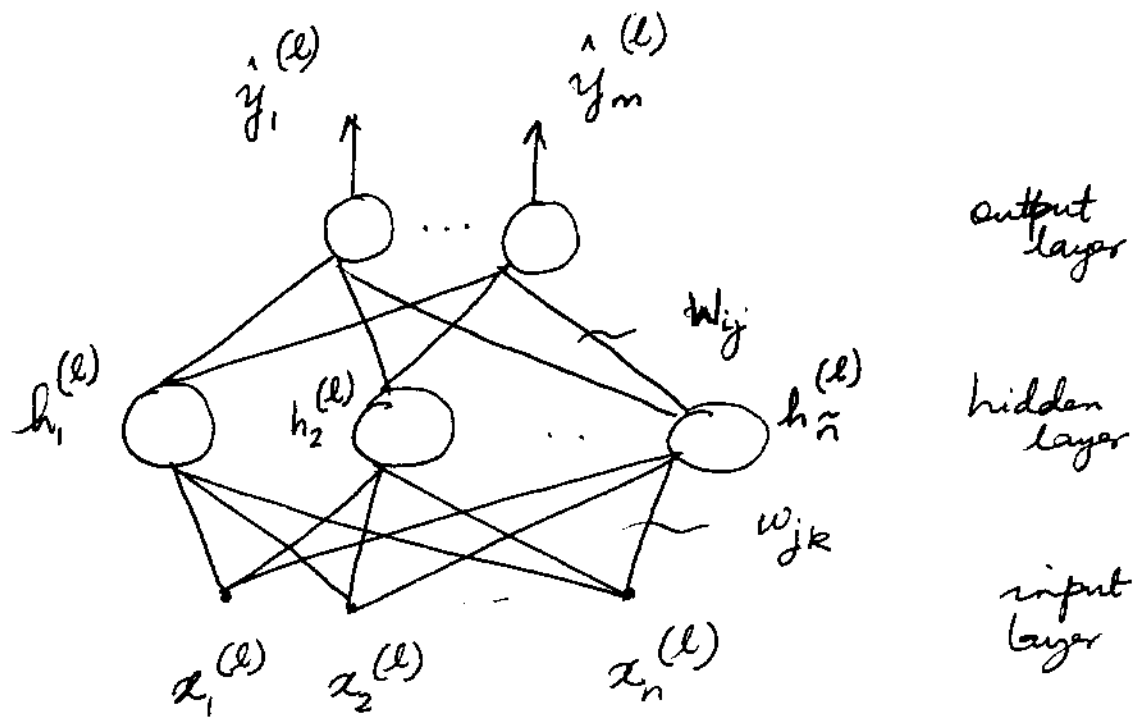


The bias is present though not shown explictly.

BP is a supervised learning algorithm.
Training data is $\left\{ \left( x^{(\ell)}, y^{(\ell)} \right) \right\}_{\ell=1}^{N}$

---

Note: we will refer to the outputs using the index $i$
       hidden layer neurons $\cdots$ $j$
       inputs. $\cdots$ $k$

so, when we say $W_{ij}$ it implies the weight between output $i$ and hid. neu. $j$

$\quad\quad$ — $\quad$ $W_{jk}$ $\cdots$ hidden neuron $j$ and input $k$.

— The hidden neuron receives as net input

$$S_j^{(\ell)} = \sum_{k=0}^{n} w_{jk} \, x_k^{(\ell)} \qquad \cdots ①$$

O/P of hidden layer neuron is :-

$$h_j^{(\ell)} = f\left(S_j^{(\ell)}\right) = f\left[\sum_{k=0}^{n} w_{jk} \, x_k^{(\ell)}\right] \qquad \cdots ②$$

— An output layer neuron receives as net input

$$S_i^{(\ell)} = \sum_{j=0}^{\hat{n}} W_{ij} \, h_j^{(\ell)}. \qquad \cdots ③$$

note: $h_0^{(\ell)} = 1$ the bias input

O/P of output layer neuron is :-

$$\hat{y}_i^{(\ell)} = f\left(S_i^{(\ell)}\right) = f\left[\sum_{j=0}^{\hat{n}} W_{ij} \, h_j^{(\ell)}\right] \qquad \cdots ④$$

Take an error (loss) function as before,

$$J^{(\ell)} = \frac{1}{2} \sum_{i=1}^{m} \left(y_i^{(\ell)} - \hat{y}_i^{(\ell)}\right)^2 \qquad \cdots ⑤$$

$$= \frac{1}{2} \sum_{i=1}^{m} e_i^{(\ell)\,2} \qquad \cdots ⑥$$

$$J = \sum_{\ell=1}^{N} J^{(\ell)}$$

Now we can use gradient descent to find
the update for the weights, i.e.

$$W_{ij} = W_{ij} - \eta \frac{\partial J^{(\ell)}}{\partial W_{ij}}$$

what is $\dfrac{\partial J^{(\ell)}}{\partial W_{ij}}$

$$\frac{\partial J^{(\ell)}}{\partial W_{ij}} = \frac{\partial J^{(\ell)}}{\partial e_i^{(\ell)}} \cdot \frac{\partial e_i^{(\ell)}}{\partial \hat{y}_i^{(\ell)}} \cdot \frac{\partial \hat{y}_i^{(\ell)}}{\partial s_i^{(\ell)}} \cdot \frac{\partial s_i^{(\ell)}}{\partial W_{ij}}$$

$$= e_i^{(\ell)} \cdot (-1) \cdot f'(s_i^{(\ell)}) \cdot h_j^{(\ell)}$$

$$= - e_i^{(\ell)} \cdot f'(s_i^{(\ell)}) \cdot h_j^{(\ell)}$$

so,

$$W_{ij} = W_{ij} + \Delta W_{ij}$$

$$= W_{ij} - \eta \frac{\partial J^{(\ell)}}{\partial W_{ij}}$$

$$= W_{ij} + \eta \, e_i^{(\ell)} \cdot f'(s_i^{(\ell)}) \cdot h_j^{(\ell)}$$

$$W_{ij} = W_{ij} + \eta \underbrace{\left[ y_i^{(\ell)} - \hat{y}_i^{(\ell)} \right] f'(s_i^{(\ell)})}_{d_i^{(\ell)}} \cdot h_j^{(\ell)} \quad \cdots \text{(7)}$$

We now know how to update the o/p to hidden weights. To get the hidden to input weights,

Recall eqn. ⑤ and ⑥

$$J^{(\ell)} = \frac{1}{2} \sum_{i=1}^{m} e_i^{(\ell)^2}$$

$$= \frac{1}{2} \sum_{i=1}^{m} \left[ y_i^{(\ell)} - \hat{y}_i^{(\ell)} \right]^2$$

$$= \frac{1}{2} \sum_{i=1}^{m} \left[ y_i^{(\ell)} - f\left[ \sum_{j=0}^{\tilde{n}} W_{ij} \cdot f\left[ \sum_{k=0}^{n} w_{jk} \cdot x_k^{(\ell)} \right] \right] \right]^2 \quad \cdots ⑧$$
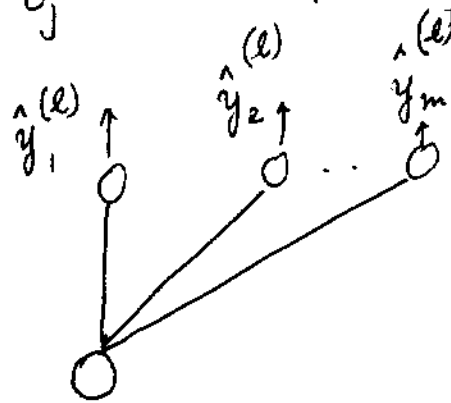
So, $\quad W_{jk} = W_{jk} - \eta \dfrac{\partial J^{(\ell)}}{\partial w_{jk}}$

what is $\dfrac{\partial J^{(\ell)}}{\partial w_{jk}}$

$$\frac{\partial J^{(\ell)}}{\partial w_{jk}} = \sum_{i=1}^{m} \frac{\partial J^{(\ell)}}{\partial e_i^{(\ell)}} \cdot \frac{\partial e_i^{(\ell)}}{\partial \hat{y}_i^{(\ell)}} \cdot \frac{\partial \hat{y}_i^{(\ell)}}{\partial s_i^{(\ell)}} \cdot \frac{\partial s_i^{(\ell)}}{\partial h_j^{(\ell)}} \cdot \frac{\partial h_j^{(\ell)}}{\partial s_j^{(\ell)}} \cdot \frac{\partial s_j^{(\ell)}}{\partial w_{jk}}$$

$$= \sum_{i=1}^{m} e_i^{(\ell)} \cdot (-1) \cdot f'\left(s_i^{(\ell)}\right) \cdot W_{ij} \cdot f'\left(s_j^{(\ell)}\right) \cdot x_k^{(\ell)}$$

$$= - \left[ \sum_{i=1}^{m} e_i^{(\ell)} \cdot f'\left(s_i^{(\ell)}\right) \cdot W_{ij} \right] f'\left(s_j^{(\ell)}\right) \cdot x_k^{(\ell)}$$

So,

$$w_{jk} = w_{jk} - \eta \frac{\partial J^{(\ell)}}{\partial w_{jk}}$$

$$= w_{jk} + \eta \left[ \sum_{i=1}^{m} e_i^{(\ell)} \cdot f'(s_i^{(\ell)}) \cdot W_{ij} \right] f'(s_j^{(\ell)}) \cdot x_k^{(\ell)}$$

$$= w_{jk} + \eta \underbrace{\left[ \sum_{i=1}^{m} \delta_i^{(\ell)} \cdot W_{ij} \right] f'(s_j^{(\ell)})}_{\delta_j^{(\ell)}} \cdot x_k^{(\ell)} \qquad \cdots \text{(9)}$$

So, $\delta_j^{(\ell)}$ is Computed as



$\delta_i^{(\ell)}$ is back-propagated through $W_{ij}$ and summed to give $\delta_j^{(\ell)}$

The name back-propagation should thus become clear.

## The Complete Algorithm

1) Initialize weights to small random values.

2) for $l = 1, 2, \ldots N$

    2a) Do the forward Pass.

$$h_j^{(l)} = f\left[s_j^{(l)}\right] = f\left[\sum_{k=0}^{n} w_{jk} x_k^{(l)}\right]$$

    and, $\hat{y}_i^{(l)} = f\left[s_i^{(l)}\right] = f\left[\sum_{j=0}^{\tilde{n}} w_{ij} h_j^{(l)}\right]$

    2b) Do the backward pass

$$\delta_i^{(l)} = \left[y_i^{(l)} - \hat{y}_i^{(l)}\right] \cdot f'\left(s_i^{(l)}\right)$$

    and,

$$\delta_j^{(l)} = \sum_{i=1}^{m} \delta_i^{(l)} \cdot w_{ij}$$
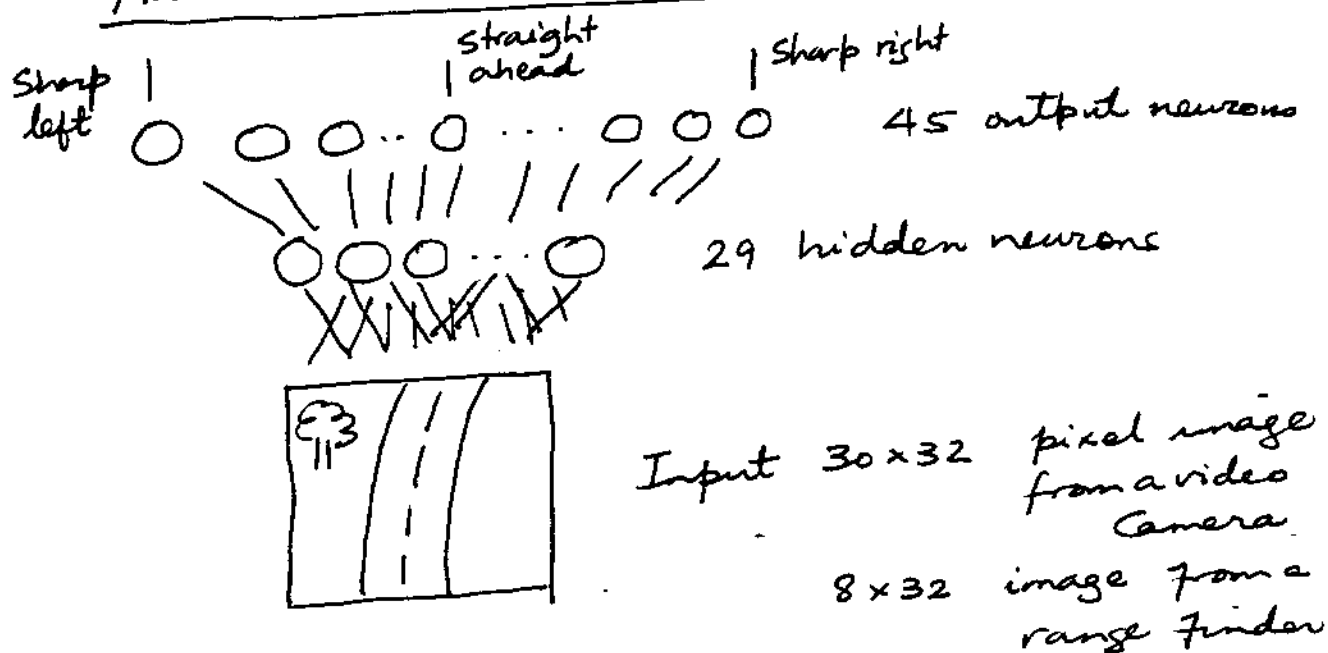
    2c) change weights

$$w_{ij} = w_{ij} + \eta\, \delta_i^{(l)}\, h_j^{(l)}$$

$$w_{jk} = w_{jk} + \eta\, \delta_j^{(\mu)} \cdot x_k^{(l)}$$

3) Present next pattern i.e. goto 2.

4) After an epoch, check J. If acceptable stop else set $l=1$ and goto step 2.

# Some Examples of Applications of Feed-Forward Networks

## Autonomous Navigation (Pomerleau)



Sharp left | straight ahead | Sharp right

45 output neurons

29 hidden neurons

Input 30 × 32 pixel image from a video camera

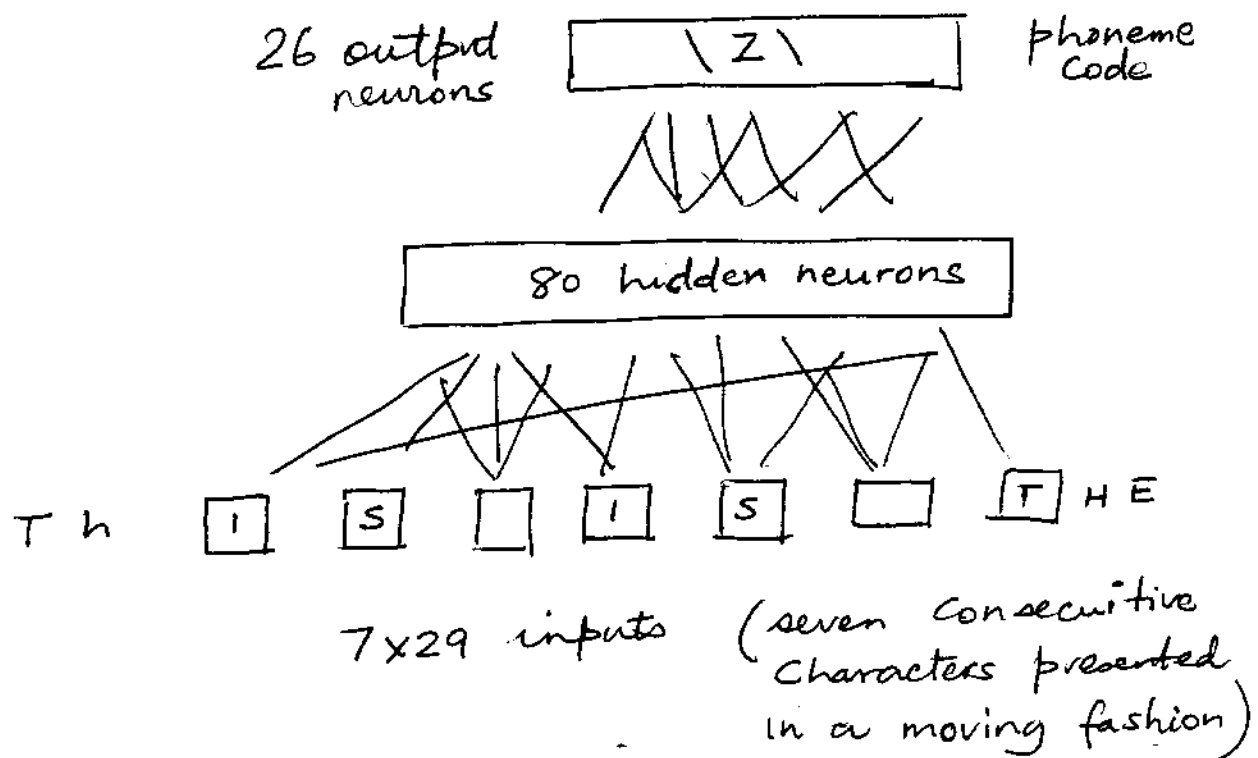8 × 32 image from a range finder

Training    1200 simulated road images
after 40 epochs, the car could autonomously
drive at 5 km/hr (3 mi/hr) on a
road through a wooded area

Speed is limited by a SUN-3 computer
to do a forward pass

Speed was twice as fast as that obtained
by any other approach.

# NETtalk (Sejnowski and Rosenberg 1987)

26 output neurons | \Z\ | phoneme code

80 hidden neurons

T h | I | S | | I | S | | T | H E

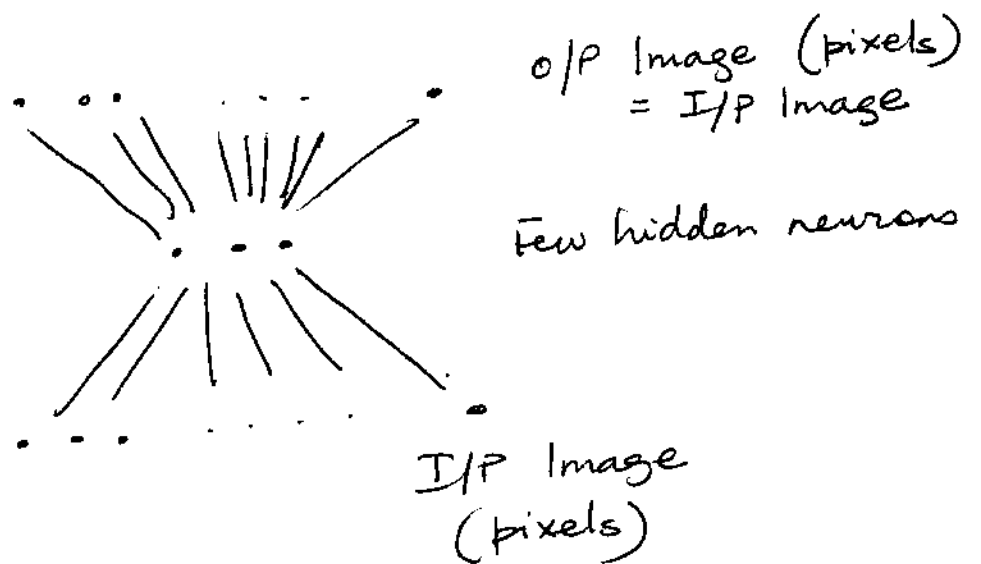7×29 inputs (seven consecutive characters presented in a moving fashion)

- Desired output was a phoneme code for the letter in the center of the window.

- Network trained using 1024 words. obtained intelligible speech after 10 training epochs (95% accuracy after 50 epochs)

- Started off like a child learning to talk and then did better.

- Some hidden neurons were found to represent meaningful properties such as the distinction between vowels and consonants.

From HKP

" It is interesting to compare NETtalk with the commercially available DECtalk which is based on hand coded linguistic rules. There is no doubt that DECtalk performs better than NETtalk but this should be seen in the context of the effort involved in creating each system. Whereas NETtalk learned from examples, the rules embodied in DECtalk are the result of about a _decade_ of analysis by _many_ _linguists_. This exemplifies the utility of Neural Networks: they are easy to construct and can be used _even_ _when_ _a_ _problem_ _is_ _not_ _fully_ _understood_. However rule based algorithms usually out perform Neural Networks _when_ _enough_ _understanding_ _is_ _available_."

# Image Compression

o/p Image (pixels)
= I/P Image

Few hidden neurons

I/P Image
(pixels)

<u>Compression Phase:</u> Train a network with patches of an image

$$o/p = I/P.$$

If the output of the hidden neurons are used, they represent a compressed code of the input (remember there are few hidden neurons)

<u>Decompression</u> use the hidden neuron o/p's and obtain the decompressed image at the o/p layer.

---

Many many other applications in Controls, Geology, Oil exploration, textiles, design, etc etc etc

# Limitations of Back-Propagation

1. Speed of training (slow)
   once trained outputs are available instantaneously.

2. Temporal Instability (moving target problem)

3. Local minima

4. Good network size (# of hidden layer neurons to use) is not known.

   However one can use pruning, growing. lateral weights for soft weight sharing.

# Advantages of FFNN's

1. General (Universal approximation)
2. Scaling - Scaling - Scaling
3. When regularization is used, provides the best overall generalization performance

# The Response of a Model

Usually, $\mathbf{y}^{(i)} = f(\mathbf{x}^{(i)}) + \epsilon$

Denote the $t$ repeated observations corresponding to this input by $\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \ldots, \mathbf{y}^{(t)}$. The optimal response $\widehat{f}^*(\mathbf{x}, \Theta)$ should minimize

$$J = \| \widehat{f}^*(\mathbf{x}; \Theta) - \mathbf{y}^{(1)} \|^2 + \| \widehat{f}^*(\mathbf{x}, \Theta) - \mathbf{y}^{(2)} \|^2 + \ldots + \| \widehat{f}^*(\mathbf{x}, \Theta) - \mathbf{y}^{(t)} \|^2$$

Minimum of $J$ is achieved when $\widehat{f}^*(\mathbf{x}, \Theta) = (\mathbf{y}^{(1)}, \mathbf{y}^{(1)}, \ldots, \mathbf{y}^{(1)})/t$, i.e. $E[\mathbf{y}|\mathbf{x}]$

# The Bias Variance Decomposition

$$
\begin{aligned}
PE &= E_{\mathcal{X}}\left[\| \left(\hat{f}(\mathbf{x}, \mathcal{X}) - E_{\mathcal{X}}\left[\hat{f}(\mathbf{x}, \mathcal{X})\right]\right) + \left(E_{\mathcal{X}}\left[\hat{f}(\mathbf{x}, \mathcal{X})\right] - E[\mathbf{y}|\mathbf{x}]\right) \|^2\right] \\
&= \underbrace{\| E_{\mathcal{X}}\left[\hat{f}(\mathbf{x}, \mathcal{X})\right] - E[\mathbf{y}|\mathbf{x}] \|^2}_{\text{Squared Bias}} + \underbrace{E_{\mathcal{X}}\left[\| \hat{f}(\mathbf{x}, \mathcal{X}) - E_{\mathcal{X}}\left[\hat{f}(\mathbf{x}, \mathcal{X})\right] \|^2\right]}_{\text{Variance}}
\end{aligned}
$$

Large k in k-NN leads to an increase
in variance, small k in large bias. This
Is typical – more complex models have
larger variance and lower bias. The
converse is true for simpler models.