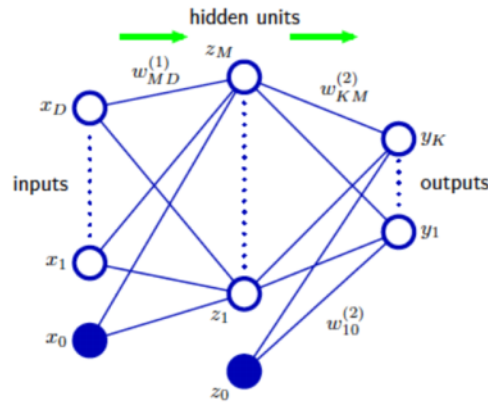


1 Backpropagation

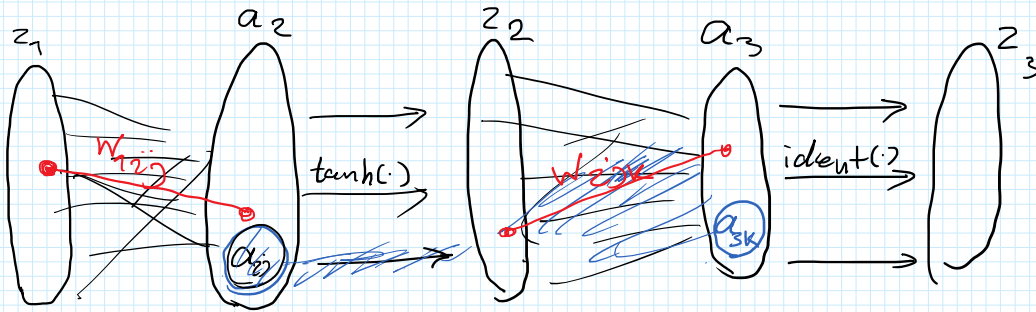
Problem 1: Apply the basic backpropagation algorithm to the network in Figure 1, with the identity $\sigma(x) = x$ as the activation function on the outputs and $h(x) = \tanh(x) = \sinh(x)/\cosh(x)$ as the activation function of the hidden neurons.



Problem 1- Backpropagation

Monday, 5 November 2018 16:01

Figure 1: Source Bishop: Figure 5.1



$$z_1^{(n)} = x_1^{(n)} \quad a_2^{(n)} = W_1^T z_1^{(n)} \quad z_2^{(n)} = \tanh(a_2^{(n)})$$

$$a_3^{(n)} = W_2^T z_2^{(n)} \quad z_3^{(n)} = a_3^{(n)}$$

$$\delta_{3k} = \frac{\partial E_n}{\partial a_{3k}} = \frac{\partial E_n}{\partial z_{3k}} \frac{\partial z_{3k}}{\partial a_{3k}} = \frac{\partial}{\partial z_{3k}} (z_{3k}^{(n)} - t_{nk})^2 \frac{1}{2} = z_{3k}^{(n)} - t_{nk}$$

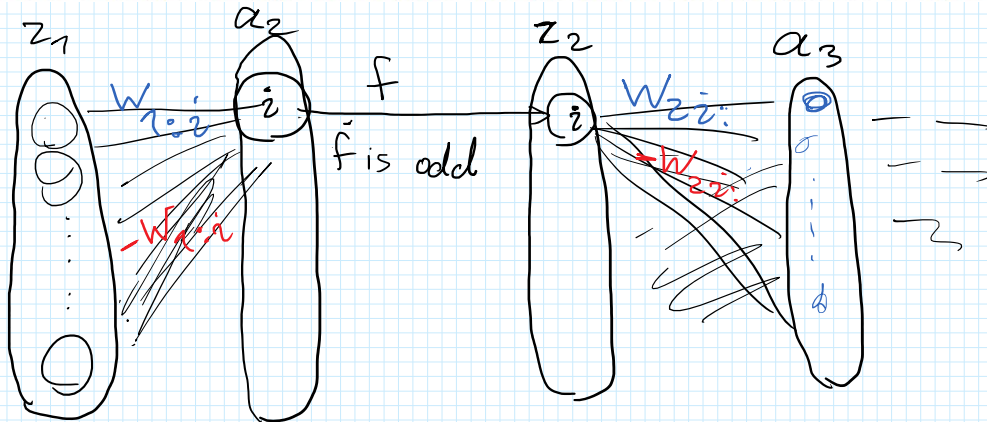
$$\delta_{2j} = \delta_3^{(n)T} W_{2j} \quad \tanh'(a_{2j}) = \delta_3^{(n)T} W_{2j} (1 - \tanh(a_{2j})^2)$$

$$\frac{\partial E_n}{\partial w_{ij}} = z_{1i}^{(n)} \delta_{2j}^{(n)}$$

$$\frac{\partial E_n}{\partial w_{2jk}} = z_{2j}^{(n)} \delta_{3k}^{(n)}$$

2 Weight Space Symmetries

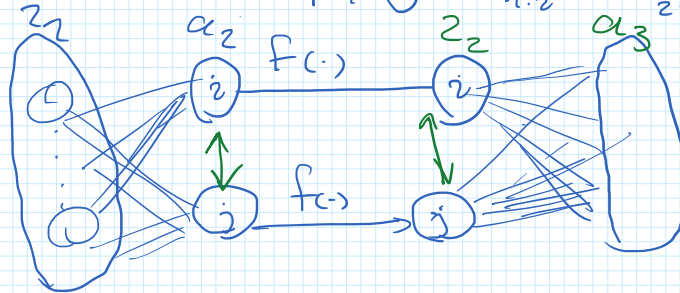
Problem 2: Assume a neural network has odd functions as non-linearities (i.e. functions for which $f(-x) = -f(x)$; e.g. \tanh). How many equivalent weight sets exist in such a neural network by considering possible changes in signs of weights and permutations among the weights?



$$a_{2i} = W_{1:i}^T z_1 \quad z_{2i} = f(w_{1:i}^T z_1) \quad a_{3j} = W_{2:i,j} f(w_{1:i}^T z_1)$$

$$a_{3j} = -w_{2ij} f(-w_{1:i}^T z_1) = w_{2ij} f(w_{1:i}^T z_1)$$

Operation 1: Multiplying $w_{1:i}$ & $w_{2:i,j}$ by -1



f has to be odd.

It is the case for e.g. $\tanh(\cdot)$.

Operation 2: Swap $[a_{2i}, z_{2i}]$ with $[a_{2j}, z_{2j}]$
(Like the swap shown by the green arrows.)

More precisely, we change the weights W to \tilde{W} as follows:

More precisely, we change the weights W to \tilde{W} as follows:

$$\tilde{W}_{1:i} = W_{1:j}, \quad \tilde{W}_{1:j} = W_{1:i}, \quad \tilde{W}_{2:i} = W_{2:j}, \quad \tilde{W}_{2:j} = W_{2:i}.$$

- About the notation: by $W_{1:j}$, we mean $\begin{bmatrix} W_{11j} \\ W_{12j} \\ \vdots \\ W_{1Dj} \end{bmatrix}$

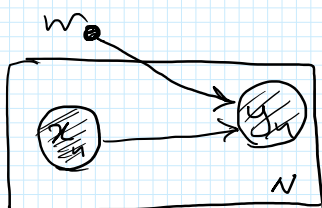
Assume we have a network like: [linear tanh] [line tanh]
with M hidden neurons.

How many equivalent networks can we get?



3 Error functions

Problem 3: Show that maximizing likelihood for a multi-class neural network model in which the network outputs that have the interpretation $f_k(\mathbf{x}, \mathbf{w}) = p(y_k = 1 | \mathbf{x})$ are represented by a softmax function is equivalent to the minimization of the cross-entropy error function. We assume that the class labels y of the training dataset $\{\mathbf{x}, y\}$ are one-hot-encoded ($y_k \in \{0, 1\}$ and $\sum_{k=0}^K y_k = 1$).

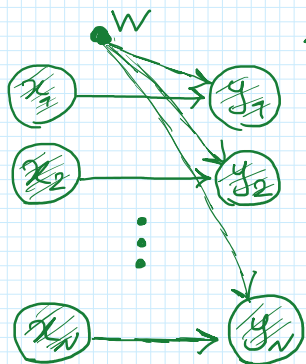


y_n : 1-of-K coding
 $y_{nk} = 1 \iff x_n \in \text{Class}_k$

$y_n \sim \text{Categorical}(y_{n1}, \dots, y_{nk} | f_1(x_n, w), \dots, f_k(x_n, w))$

Side note about this figure:

This graphical model is equivalent to this figure:



But using the above form
 (known as plate notation)
 is more convenient and more common.

Learning w using MLE:

$$P(\text{All}) = \prod_n P(y_n | x_n) \left(\prod_n P(x_n) \right)$$

$$= \prod_n \text{Categorical}(y_{n1}, \dots, y_{nk} | f_1(x_n, w), \dots, f_k(x_n, w))$$

$$= \prod_n \prod_k f_k(x_n, w)^{y_{nk}}$$

$$\Rightarrow \log \text{likelihood} = \sum_n \sum_k y_{nk} \log f_k(x_n; w) = (-1) \text{Cross-Entropy}$$

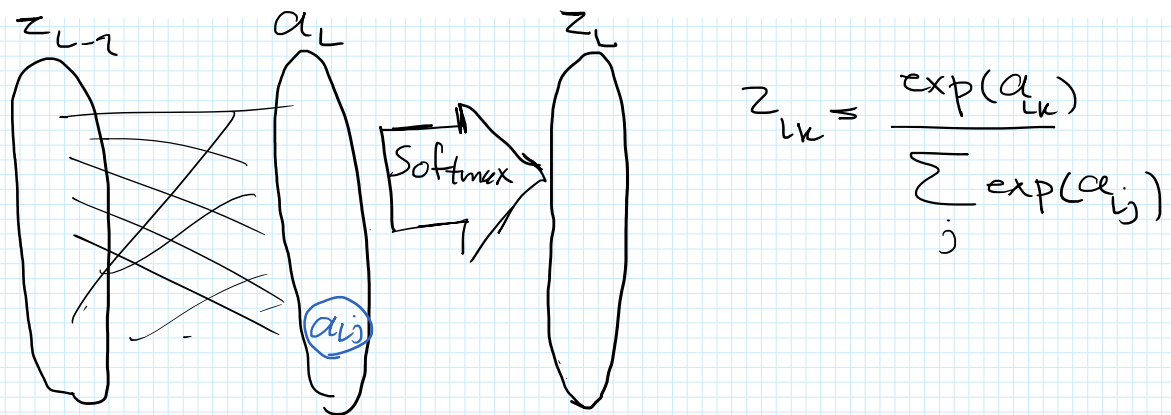
$$w_* = \underset{w}{\operatorname{argmax}} \log \text{likelihood} = \underset{w}{\operatorname{argmin}} [\text{Cross-Entropy}]$$

Problem 4: Show that the derivative of the standard (multi-class) cross-entropy error function

$$E(w) = \sum_{n=1}^N E_n(w) = - \sum_{n=1}^N \sum_{k=1}^K y_k^{(n)} \log f_k(\mathbf{x}^{(n)}, \mathbf{w})$$

with respect to the activation a_k for the output units with a softmax activation function satisfies

$$\frac{\partial E_n}{\partial a_k} = f_k(\mathbf{x}^{(n)}, \mathbf{w}) - y_k^{(n)}$$



Recall: Chain-rule in higher dim space.

$$\frac{\partial E_n}{\partial a_L} = \frac{\partial E_n}{\partial z_L^{(n)}} \frac{\partial z_L^{(n)}}{\partial a_L}$$

\downarrow \downarrow \downarrow
 $1 \times K$ $1 \times K$ $K \times K$

$$\begin{matrix} A & B & C \\ \mathbb{R}^I & \mathbb{R}^J & \mathbb{R}^K \end{matrix}$$

$$\frac{\partial C}{\partial A} = \frac{\partial C}{\partial B} \frac{\partial B}{\partial A}$$

\downarrow \downarrow \downarrow
 $K \times I$ $K \times J$ $J \times I$

$$\Rightarrow \frac{\partial E_n}{\partial a_{ij}} = \frac{\partial E_n}{\partial z_L^{(n)}} \frac{\partial z_L^{(n)}}{\partial a_{ij}}$$

$$\frac{\partial E_n}{\partial z_{Li}^{(n)}} = \frac{\partial}{\partial z_{Li}^{(n)}} \left[\sum_{k=1}^K y_{nk} \log(z_{Lk}^{(n)}) \right] = - \frac{y_{ni}}{z_{Li}^{(n)}}$$

$$\frac{\partial z_{Li}^{(n)}}{\partial a_{ij}} = \frac{\partial}{\partial a_{ij}} \frac{\exp(a_{Li})}{\sum_j \exp(a_{Lj})}$$

$$i=j \Rightarrow \frac{\partial z_{Li}^{(n)}}{\partial a_{ii}} = \frac{\exp(a_{Li}) \left[\sum_j \exp(a_{Lj}) - \exp(a_{Li}) \right]}{\left(\sum_j \exp(a_{Lj}) \right)^2}$$

$$i=j \Rightarrow \frac{\partial z_{li}^{(n)}}{\partial a_{li}} = \frac{\exp(a_{li}) [\sum_k \exp(a_{lk}) - \exp(a_{li})]}{[\sum_k \exp(a_{lk})]^2}$$

$$= \cancel{z_{li}^{(n)} - z_{li}^{(n)}}$$

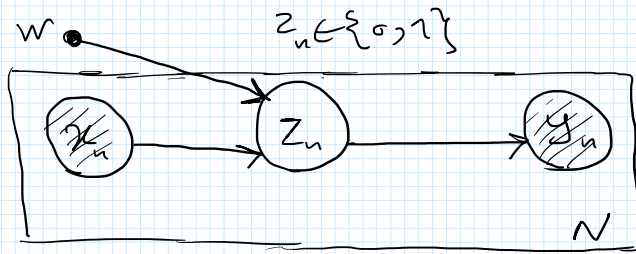
$$i \neq j \Rightarrow \frac{\partial z_{li}^{(n)}}{\partial a_{lj}} = \frac{-\exp(a_{li}) \exp(a_{lj})}{[\sum_k \exp(a_{lk})]^2} = \cancel{-z_{li}^{(n)} z_{lj}^{(n)}}$$

$$\Rightarrow \frac{\partial E_n}{\partial a_{lj}} = \sum_i \left(-\frac{y_{ni}}{\cancel{z_{li}^{(n)}}} \right) \left(\cancel{-z_{li}^{(n)}} z_{ij}^{(n)} \right) - \cancel{z_{lj}^{(n)}} \frac{y_{nj}}{\cancel{z_{lj}^{(n)}}}$$

$$= z_{lj}^{(n)} \left(\sum_i \overset{1}{\cancel{+} y_{ni}} \right) - y_{nj} = z_{lj}^{(n)} - y_{nj} = \delta_{lj}^{(n)}$$

4 Robust classification

Problem 5: Consider a binary classification problem in which the target values are $y \in \{0, 1\}$, with a network output $f(\mathbf{x}, \mathbf{w})$ that represents $p(y = 1 | \mathbf{x}, \mathbf{w})$, and suppose that there is a probability ε that the class label on a training data point has been incorrectly set. Assuming independent and identically distributed data, write down the error function corresponding to the negative log likelihood. Verify that the well known error function for binary classification is obtained when $\varepsilon = 0$. Note that this error function makes the model robust to incorrectly labelled data, in contrast to the usual error function.



likelihood is a function of z_n

$$P(z_n | x_n) = \text{Bernoulli}(f(x_n, w))$$

$$P(y_n | z_n) = \begin{array}{c|cc} z_n & P(y_n=0|z_n) & P(y_n=1|z_n) \\ \hline 0 & 1-\varepsilon & \varepsilon \\ 1 & \varepsilon & 1-\varepsilon \end{array}$$

To maximize the Likelihood, we need to maximize:

$$\prod_n P(x_n, z_n, y_n)$$

But the problem is that z_n is not observed.

Thus we can maximize: $\prod_n P(x_n, y_n)$

$$\text{Where } p(x_n, y_n) = p(x_n, \underline{z_n=0}, y_n) + p(x_n, \underline{z_n=1}, y_n)$$

For notational convenience, we write it as:

$$p(x_n, y_n) = p(x_n, \underline{z_n=y_n}, y_n) + p(x_n, \underline{z_n=1-y_n}, y_n)$$

\Rightarrow

$$P(A11) = \prod_n \left[p(x_n, \underline{z_n=y_n}, y_n) + p(x_n, \underline{z_n=1-y_n}, y_n) \right]$$

$$= \prod_n \left[\underbrace{p(x_n)}_{\text{X}} \underbrace{\text{Bernoulli}(y_n, f(x_n, w))}_{\text{X}} (1-\varepsilon) \right]$$

Recall:

$$\text{Bernoulli}(x, \theta) = \theta^x (1-\theta)^{1-x}$$

$$= \prod_n \left[f(x_n; w)^{y_n} (1 - f(x_n; w))^{1-y_n} \right] \\ + \underbrace{p(x_n)}_{\text{Bernoulli}(1-y_n, f(x_n; w))} \xi$$

recall:

$$\text{Bern}(x; \theta) = \theta^x (1-\theta)^{1-x}$$

$$= \prod_n \left[\underbrace{f(x_n; w)^{y_n} (1 - f(x_n; w))^{1-y_n}}_{\text{Bernoulli}(y_n, f(x_n; w))} (1 - \xi) \right] \\ + \underbrace{f(x_n; w)^{1-y_n} (1 - f(x_n; w))^{y_n}}_{\text{Bernoulli}(1-y_n, f(x_n; w))} (\xi)$$

In the case $\xi = 0 \Rightarrow \text{Cost} = - \prod_n \left[f(x_n; w)^{y_n} (1 - f(x_n; w))^{1-y_n} \right]$

$$\Rightarrow \log \text{Cost} = - \sum_n y_n \log f(x_n; w) + (1-y_n) \log (1 - f(x_n; w))$$

which is the binary cross-entropy cost function.