# Task 1.1

Over a single data point we get:

$$-E(w) = t^n \ln(y^n) + (1-t^n) \ln(1-y^n)$$

using $y^n = g_w^n$ and $\dfrac{\partial g_w^n}{\partial w_j} = x_j^n g_w^n (1-g_w^n)$

allows the usage of the chain rule to find $\dfrac{\partial -E^n(w)}{\partial w_j}$

$$\frac{\partial (t^n \ln(g_w^n))}{\partial w_j} = \frac{t^n \cancel{(x_j^n \; g_w^n} (1-g_w^n)}{\cancel{g_w^n}} \rightarrow \text{example chain rule}$$

$$\frac{\partial ((1-t^n) \ln(1-g_w^n))}{\partial w_j} = -\frac{(1-t^n)}{\cancel{1-g_w^n}} (x^n g_w^n) \cancel{(1-g_w^n)}$$

$$\Rightarrow \quad \frac{-\partial E^n(w)}{\partial w_j} = t^n x_j^n (1-g_w^n) - (1-t^n)(x_j^n g_w^n)$$

$$= t^n x_j^n - t^n x_j^n g_w^n - (x_j^n g_w^n - t^n x_j^n g_w^n)$$

$$= t^n x_j^n \cancel{- t^n x_j^n g_w^n} - x_j^n g_w^n \cancel{+ t^n x_j^n g_w^n}$$

$$= t^n x_j^n - x_j^n g_w^n = (t^n - g_w^n) x_j^n$$

$$= \underline{(t^n - y^n) x_j^n} \quad \blacksquare$$

$$y_k^n = \frac{e^{w_k^T x^n}}{\sum_{k'} e^{w_{k'}^T x^n}}$$

Using the quotient rule:

$$\frac{\partial y_k^n}{\partial w_j} = \frac{\frac{\partial}{\partial w_j} e^{w_k^T x^n}}{\left(\sum_{k'} e^{w_{k'}^T x^n}\right)^2} \cdot \sum_{k'} e^{w_{k'}^T x^n} - \frac{e^{w_k^T x^n}}{\left(\sum_{k'} e^{w_{k'}^T x^n}\right)^2} \cdot \frac{\partial}{\partial w_j} \sum e^{w_k^T x^n}$$

Now: we split this into parts of $k=j$ and $k \neq j$, such that we differ between whether the weights are inputs to the wanted output or not:

$\boxed{k=j:}$ => $\frac{x^n e^{w_k^T x^n}}{\sum_{k'} e^{w_{k'}^T x^n}} - \frac{e^{w_k^T x^n}}{\left(\sum_{k'} e^{w_{k'}^T x^n}\right)^2} \cdot x^n e^{w_{\textcircled{j}}^T x^n}$  ← mind this $j$!

=> $\frac{\partial y_k^n}{\partial w_j} = x^n y_k^n - x^n y_k^n y_j^n = \underline{x^n y_k^n (1 - y_j^n)}$

$\boxed{k \neq j:}$  $0 - x^n y_k^n y_j^n = \frac{\partial y_k^n}{\partial w_j}$  (by the same method as above for $k=j$)

With the derivative of softmax, we may obtain the gradient of the loss function:

$$\frac{\partial E}{\partial w_{kj}} = -\sum_k^C t_k \ln(y_k^n) = -\sum_{k=1}^C \frac{t_k}{y_k^n} \frac{\partial y_k^n}{\partial w_{kj}}$$  (move the minus sign)

$$\frac{\partial y_k^n}{\partial w_{kj}} = \begin{cases} x^n y_k^n (1 - y_j^n), & k=j \\ -x^n y_k^n y_j^n, & k \neq j \end{cases} => y_k^n x^n (\underset{\text{Kronecker delta}}{\delta_{kj}} - y_{\textcircled{j}}^n)$$

Thus: $-\frac{\partial E}{\partial w_{kj}} = \sum_{k=1}^C \frac{t_k}{y_k^n} y_k^n (\delta_{kj} - y_k^n)$

$$= \sum_{k=1}^C t_k x^n \delta_{ij} - \sum_{k=1}^C t_k x^n y_j^n = x^n (t_k - y_k^n) \blacksquare$$

(if we are using the weight $w_{kj}$, making $j=k$ here, not to be confused with the $j=k$ from the derivation)

# Task 2.2 a)

$$\mathcal{J}(w) = \mathcal{E}(w) + \lambda C(w)$$

$$\frac{\mathcal{J}(w)}{\partial w} = \boxed{\frac{\mathcal{E}(w)}{\partial (w)}} + \boxed{\lambda \frac{C(w)}{\partial w}}$$

↳ given from previous task

$$\frac{\partial C(w)}{\partial w} = \frac{\partial \left( \lambda \sum_{i,j} w_{i,j}^2 \right)}{\partial w} = \lambda 2 w \quad \boxed{}$$

( The summation disappears, as we are differentiating over one specific weight, whilst the others remain constant).

I've also seen $\frac{\lambda}{2n} C(w)$ used, which gives:

$$\frac{\partial \frac{\lambda}{2n} \mathcal{E} w^2}{\partial w} = \underline{\frac{\lambda}{n} w}.$$