

Question 4(e):

Let  $m$  denote the number of features

$$\begin{aligned}
 \frac{\partial J}{\partial w} * N &= X^T (y - t) \\
 &= \begin{matrix} x_1^{(1)} & \dots & x_1^{(N)} & y^{(1)} - t^{(1)} \\ \dots & \dots & \dots & \dots \\ x_m^{(1)} & \dots & x_m^{(N)} & y^{(N)} - t^{(N)} \end{matrix} * \\
 \frac{\partial J}{\partial w} &= \left(\frac{1}{N}\right) * \begin{matrix} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_1^{(i)} \\ \dots \\ \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_{1m}^{(i)} \end{matrix} \\
 &= \begin{matrix} \partial J / \partial w_1 \\ \dots \\ \partial J / \partial w_m \end{matrix} \\
 \therefore \left[\frac{X^T (y - t)}{N}\right]_j &= \left(\frac{1}{N}\right) \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_1^{(i)} \\
 \therefore \left[\frac{\partial J}{\partial w}\right]_j &= \left[\frac{X^T (y - t)}{N}\right]_j \leftrightarrow \frac{\partial J}{\partial w_j} = \left(\frac{1}{N}\right) \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_1^{(i)}
 \end{aligned}$$

Question 4(f):

**Prove:**  $L_{LCE}(z, t) = L_{CE}(\text{sigmoid}(z), t)$

$$L_{CE}(y, t) = -t * \log(y) - (1 - t) * \log(1 - y)$$

**Substitute  $y$  for sigmoid( $z$ )**

$$\begin{aligned}
 L_{ce}(\text{sigmoid}(z), t) &= -t * \log\left(\frac{1}{1 + e^{-z}}\right) - (1 - t) * \log\left(1 - \frac{1}{1 + e^{-z}}\right) \\
 &= t * \log(1 + e^{-z}) - (1 - t) * \log\left(\frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}}\right) \\
 &= t * \log(1 + e^{-z}) - (1 - t) * \log\left(\frac{e^{-z}}{1 + e^{-z}}\right) \\
 &= t * \log(1 + e^{-z}) + (1 - t) * \log\left(\frac{1 + e^{-z}}{e^{-z}}\right) \\
 &= t * \log(1 + e^{-z}) + (1 - t) * \log(1 + (e^{-z})^{-1})
 \end{aligned}$$

$$= t * \log(1 + e^{-z}) + (1 - t) * \log(1 + e^z)$$

$$\therefore L_{LCE}(z, t) = L_{CE}(\text{sigmoid}(z), t)$$

Question 6(e):

I suspect that the biggest reason behind both occurrences is the fact that there are more differences between 4 and 7 than between 5 and 6. A 5 is more likely to look like a 6 and vice versa than a 4 is to a 7.

Question 6(f):

Because we only have two classes, an odd number of  $k$  prevents a tie in the decision. In each prediction, one class will hold the majority of predictions

Question 6(g):

Because numbers were designed to look different from one another and are unlikely to have many outliers in comparison to some other data sets. This lack of outliers makes it a great candidate for KNN because any given hand drawn digit is highly likely to look like another hand drawn version of the same digit