

② c) $w_{NN}^* \geq 0$

$$\mathcal{L}_{NN}(w_{NN}^*) = \frac{1}{L} \sum_{n=1}^N (f(w_{NN}^*(x_n)) - y_n)^2$$

$$\mathcal{L}_{LS}(w_{LS}^*) = \frac{1}{L} \sum_{n=1}^N (w_{LS}^{*T} x_n - y_n)^2$$

$$\begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \dots \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = 0$$

since $x^0 \geq 0$

$$w_1^* \dots w_{L-1}^* \geq 0$$

$$\Rightarrow \text{relu}(w; l) \geq 0, \text{ with } l = \text{relu}(w_1 \dots w_{i-1} x)$$

inputs
disjunctive w/

$$x^0 \quad (0,1) \quad (0,1) \quad (0,1) \quad \hat{}$$

L layers with relu $\rightarrow \min(0, x)$

$$\text{relu} = \min(0, x) \text{ but } x \geq 0 \Rightarrow f(x) = \text{relu}(x) = x$$

$$\Rightarrow f(w_{NN}^*(X)) = w_{L+1}^{*T} \left(\dots \text{relu} \left(w_2^{*T} \left(\text{relu} \left(w_1^{*T} X \right) \right) \right) \right) = w_{L+1}^{*T} \dots w_1^{*T} X$$

$$\Rightarrow \mathcal{L}_{NN} = \frac{1}{L} \sum_{n=1}^N \left[\underbrace{\left(w_{L+1}^{*T} \dots w_1^{*T} X - y \right)^2}_{w_{NN}^*} \right]_n$$

$$\cdot \mathcal{L}_{LS} = \frac{1}{L} \sum_{n=1}^N (w_{LS}^{*T} x_n - y_n)^2 =$$

successive multiplication of matrices has no further "resolution" (same matrix dimension) than a singular matrix, therefore the optimum solution to the NN loss function would yield the same weight matrix as the standard LS loss function. Thus, $L_{NN} = L_{LS}$

b) in this case: $\mathcal{L}_{NN} = \frac{1}{L} \sum_{n=1}^N \left[\left(w_{L+1}^{*T} \left(\dots \text{relu} \left(w_2^{*T} \left(\text{relu} \left(w_1^{*T} X \right) \right) \right) \right) - y \right)^2 \right]_n$

Considering that the w^* for the least squares loss returned a non-negative weight vector, combined with the samples being non-negative, we should assume that the regression targets are also non-negative. . . .

With this in mind, and considering that the feed forward neural network can still attain negative weights in the output and hidden layers, we should expect the feed forward neural network to have more accurate predictions. Moreover, since every hidden layer is coupled with an activation function and, therefore, should be able to extract features not exposed in the raw samples, the only situation where any particular weight would be negative, should be when a extracted feature is detrimental to the prediction. Thus, since the LS loss does not have access to this specific feature it can't compensate resulting in a higher combined loss over all samples.