

Problem 2: You are trying to solve a regression task and you want to choose between two approaches:

1. A simple linear regression model.
2. A feed forward neural network $f_{\mathbf{W}}(\mathbf{x})$ with L hidden layers, where each hidden layer $l \in \{1, \dots, L\}$ has a weight matrix $\mathbf{W}_l \in \mathbb{R}^{D \times D}$ and a ReLU activation function. The output layer has a weight matrix $\mathbf{W}_{L+1} \in \mathbb{R}^{D \times 1}$ and no activation function.

In both models, there are no bias terms.

Your dataset \mathcal{D} contains data points with nonnegative features \mathbf{x}_n and the target y_n is continuous:

$$\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N, \quad \mathbf{x}_n \in \mathbb{R}_{\geq 0}^D, \quad y_n \in \mathbb{R}$$

Let $\mathbf{w}_{LS}^* \in \mathbb{R}^D$ be the optimal weights for the linear regression model corresponding to a *global* minimum of the following least squares optimization problem:

$$\mathbf{w}_{LS}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^D} \mathcal{L}_{LS}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^D} \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2$$

Let $\mathbf{W}_{NN}^* = \{\mathbf{W}_1^*, \dots, \mathbf{W}_{L+1}^*\}$ be the optimal weights for the neural network corresponding to a *global* minimum of the following optimization problem:

$$\mathbf{W}_{NN}^* = \arg \min_{\mathbf{W}} \mathcal{L}_{NN}(\mathbf{W}) = \arg \min_{\mathbf{W}} \frac{1}{2} \sum_{n=1}^N (f_{\mathbf{W}}(\mathbf{x}_n) - y_n)^2$$

- a) Assume that the optimal \mathbf{W}_{NN}^* you obtain are non-negative.

What will the relation ($<, \leq, =, \geq, >$) between the neural network loss $\mathcal{L}_{NN}(\mathbf{W}_{NN}^*)$ and the linear regression loss $\mathcal{L}_{LS}(\mathbf{w}_{LS}^*)$ be? Provide a mathematical argument to justify your answer.

Note that for any non-negative x and non-negative W it holds $\text{ReLU}(xW) = xW$

Therefore, since our data points have non-negative features \mathbf{x}_i and the optimal weights \mathbf{W}_{NN}^* are non-negative, every ReLU layer is equivalent to a linear layer when plugging in the optimal weights. This means we can write:

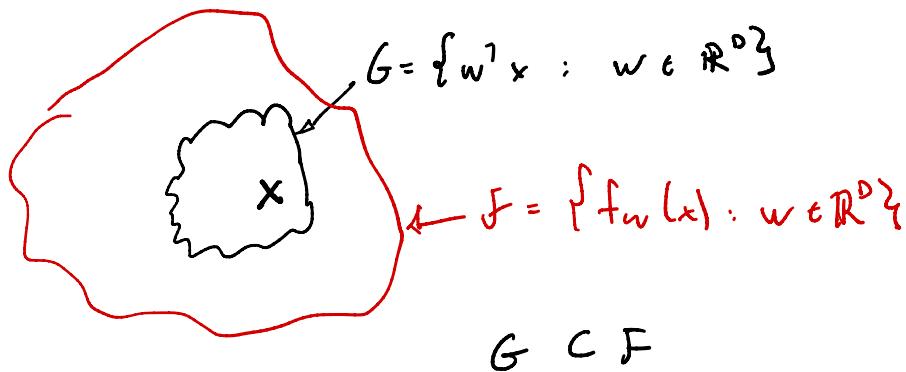
$$\begin{aligned} f_{\mathbf{W}_{NN}^*}(\mathbf{x}_i) &= \text{ReLU}(\text{ReLU}(\text{ReLU}(\mathbf{x}_i^T \mathbf{W}_1^*) \mathbf{W}_2^*) \dots \mathbf{W}_L^*) \mathbf{W}_{L+1}^* \\ &= \mathbf{x}_i^T \mathbf{W}_1^* \mathbf{W}_2^* \dots \mathbf{W}_{L+1}^* \\ &= \mathbf{x}_i^T \mathbf{W}_{NN}^* \end{aligned}$$

where we defined $\mathbf{W}_{NN}^* = \mathbf{W}_1^* \mathbf{W}_2^* \dots \mathbf{W}_{L+1}^*$. From this we can see that the neural network with optimal weights behaves like a linear regression with a different set of weights \mathbf{W}_{NN}^* . Note also that linear regression is a special case of the above neural network, i.e. for any weights \mathbf{W}_{LS} you can find weights \mathbf{W}_{NN} that produces the same output.

Given the above facts and since the optimal weights correspond to a global minimum, we can conclude that $\mathcal{L}_{NN}(W_{NN}^*) = \mathcal{L}_{LS}(w_{LS}^*)$ and the optimal weights found by solving the least squares optimization problem will be $w_{LS}^* = W_{NN}^*$.

*

Not every NN can be represented with LR model (putting negative weights)



$$\min_{f \in \mathcal{F}} \frac{1}{2} \sum_{n=1}^N (f(x_n) - y_n)^2 \geq \min_{f \in \mathcal{G}} \frac{1}{2} \sum_{n=1}^N (f(x_n) - y_n)^2$$

b) In contrast to (a), now assume that the optimal weights w_{LS}^* you obtain are non-negative.

What will the relation ($<$, \leq , $=$, \geq , $>$) between the linear regression loss $\mathcal{L}_{LS}(w_{LS}^*)$ and the neural network loss $\mathcal{L}_{NN}(W_{NN}^*)$ be? Provide a mathematical argument to justify your answer.

As stated in (a) linear regression is a special case of the above neural network, i.e. for any weights w_{LS} you can find weights W_{NN} that produces the same output. That is, everything can be learned with a linear regression can be learned equally well with neural networks. However, the reverse direction doesn't hold, since in principle neural networks can learn more complicated functions compared to linear regression. Moreover, the given fact that w_{LS}^* are non-negative does not tell us anything about the optimal weights of the neural network W_{NN}^* . Therefore it holds $\mathcal{L}_{NN}(W_{NN}^*) \leq \mathcal{L}_{LS}(w_{LS}^*)$ since the neural network can potentially find a better fit for the data (e.g. by taking advantage of non-linearity).

