

Homework

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 $f_W(x)$ with L hidden layers, where each hidden layer $l \in \{1, \dots, L\}$ has a weight matrix $W_l \in \mathbb{R}^{D_{l+1} \times D_l}$ and a ReLU activation function. The output layer has a weight matrix $W_{L+1} \in \mathbb{R}^{D \times 1}$ and no activation function.

In both models, there are no bias terms.

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

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

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

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

Let $W_{NN}^* = \{W_1^*, \dots, W_{L+1}^*\}$ be the optimal weights for the neural network corresponding to a *global* minimum of the following optimization problem:

$$W_{NN}^* = \arg \min_W \mathcal{L}_{NN}(W) = \arg \min_W \frac{1}{2} \sum_{n=1}^N (f_W(x_n) - y_n)^2$$

a) Assume that the optimal W_{NN}^* you obtain are non-negative.

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

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

Therefore, since our data points have non-negative features x_i and the optimal weights 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

$$\begin{aligned} f_{W_{NN}^*}(x_i) &= \text{ReLU}(\text{ReLU}(\text{ReLU}(x_i^T W_1^*) W_2^*) \dots W_L^*) W_{L+1}^* \\ &= x_i^T W_1^* W_2^* \dots W_{L+1}^* \\ &= x_i^T w_{NN}^* \end{aligned}$$

where we defined $w_{NN}^* = W_1^* W_2^* \dots W_{L+1}^*$. From this we can see that the neural network with optimal weights behaves like a linear regression with a different set of weights w_{NN}^* .

Note also that 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 produce the same output.

Given the above facts and since we the optimal weights correspond to a global minima 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 $w_{LS}^* = w_{NN}^*$.

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 produce the same output. That is, everything that can be learned with a linear regression can be learned equally well with a neural network.

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}^* .

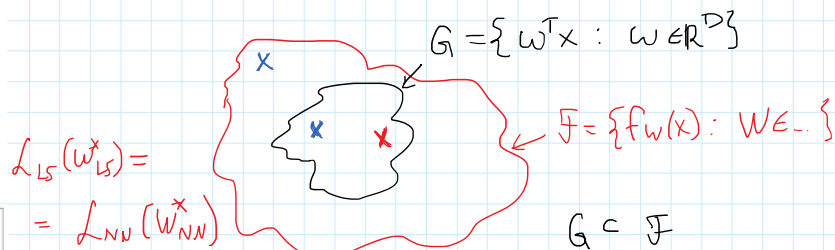
Therefore it holds $\mathcal{L}_{NN}(W_{NN}^*) \leq \mathcal{L}_{LS}(w_{LS}^*)$ since the neural network can potentially find a better fit for the data (e.g. by taking advantage of non-linearity).

Any LR model $g(x) = w^T x$ I can represent with a NN model $f_W(x)$.

$$f_W(x) = \text{ReLU}(\text{ReLU}(x^T I) I) \cdot w = x^T w = w^T x$$

$$\text{ReLU}(Wx) = Wx \quad \text{for } W \geq 0, x \geq 0$$

Not every NN I can represent with a LR model.



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

If $W_{NN}^* \geq 0$, then $f_W(x) = w^T x$

$$L_{LS}(w_{LS}^*) \geq L_{NN}(W_{NN}^*)$$

$$\begin{matrix} \xrightarrow{1} \\ 0 \end{matrix} \begin{matrix} \xrightarrow{M} \\ 1 \end{matrix} = x$$

$$\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_N \quad x \cdot \text{mean}(\text{dim}=1)$$

$$\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_N \quad x \cdot \text{mean}(\text{dim}=1, \text{keepdim=True})$$

$$\frac{d}{dx} \frac{1}{\sqrt{x^2 + \epsilon}}$$

$$x \quad [N, C, H, W]$$

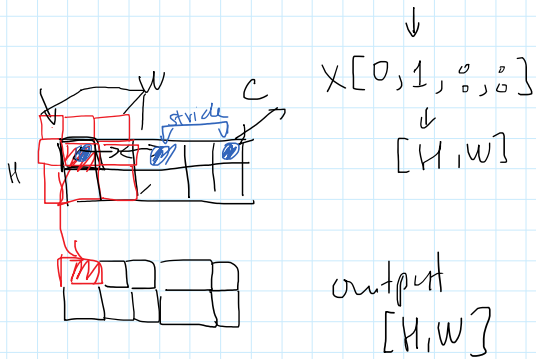
$$\mathcal{R}(x) = \mathcal{F}(x) + x$$

in_channels

out_channels

kernel_size = 3

stride = 1



$[N, C, H, W]$

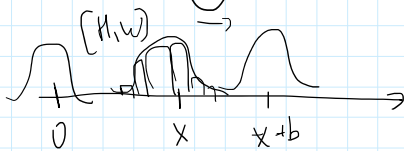
C_{in} } provided by user

C_{out} }

$k = 3$

padding = 1

stride = $\begin{cases} 1 \Rightarrow \text{out } [N, C, H, W] \\ 2 \Rightarrow \text{out } [N, C, \frac{H}{2}, \frac{W}{2}] \end{cases}$



$[N, C, H, W] \leftarrow x$

$[N, C, \frac{H}{2}, \frac{W}{2}] \leftarrow \mathcal{F}(x)$



output shape

$C \quad H \quad W$
 $[N, 3, 32, 32]$

conv ($k=3, s=1$) $[N, 16, 32, 32]$

stack ($C_{in}=16, k=3, C_{out}=16, S=1$) $[N, 16, 32, 32]$

stack ($C_{in}=16, k=3, C_{out}=32, S=2$) $[N, 32, 16, 16]$

stack ($C_{in}=32, k=3, C_{out}=64, S=2$) $[N, 64, 8, 8]$

def stack (n):
 { Residual block (-)
 ...
 Residual block (-)

ResBlock: $\left. \begin{array}{l} (s=2) \text{ conv 1} \\ \text{conv 2} \end{array} \right\} 2 \text{ layers}$

Res Stack ($s=2$) $\left. \begin{array}{l} \text{ResBlock 1} \\ \vdots \\ \text{ResBlock n} \end{array} \right\} \begin{array}{l} 2n \text{ layers} \\ = \\ n \cdot \text{ResBlocks} \end{array}$ c h w

x $[N, 3, 32, 32]$

conv ($k=3, s=1$) $[N, 16, 32, 32]$

stack ($C_{in}=16, k=3, C_{out}=16, S=1$) $[N, 16, 32, 32]$

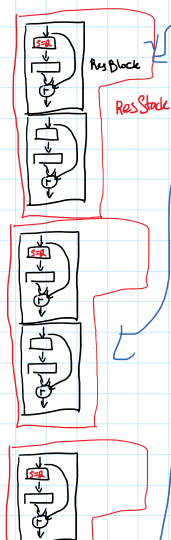
stack ($C_{in}=16, k=3, C_{out}=32, S=2$) $[N, 32, 16, 16]$

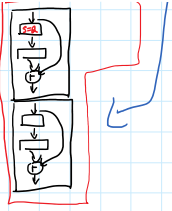
stack ($C_{in}=32, k=3, C_{out}=64, S=2$) $[N, 64, 8, 8]$

ArgPool $[N, 64, 1]$

squeeze $[N, 64]$

Linear $[N, 10]$





Linear $[N, 10]$

1 conv layer

6n conv. layers

1 fully connected

