

A Proofs of lower bounds

A.1 Softmax output layers

Let f be a linear classifier. For an observation \mathbf{x} and any class $c \neq f(\mathbf{x})$, we can compute the norm of the smallest perturbation \mathbf{r} such that $f(\mathbf{x} + \mathbf{r}) = c$. Let us denote this perturbation by $\Delta_{\text{adv}}(\mathbf{x}; f, c)$:

$$\Delta_{\text{adv}}(\mathbf{x}; f, c) = \min_{\mathbf{r} \in \mathbb{R}^d} \{\|\mathbf{r}\| \mid f(\mathbf{x} + \mathbf{r}) = c\}. \quad (5)$$

It is easily seen that

$$\Delta_{\text{adv}}(\mathbf{x}; f) = \min_{c \neq f(\mathbf{x})} \{\Delta_{\text{adv}}(\mathbf{x}; f, c)\}. \quad (6)$$

The restriction that $c \neq f(\mathbf{x})$ is important, since otherwise we would get a degenerate solution of $\Delta_{\text{adv}}(\mathbf{x}; f) = 0$ for all \mathbf{x} . Intuitively, $\Delta_{\text{adv}}(\mathbf{x}; f)$ is the norm of the projection of \mathbf{x} onto the class “closest” to \mathbf{x} but distinct from $f(\mathbf{x})$. A necessary condition for any adversarial perturbation \mathbf{q} is given by the following

Lemma A.1. *Let f be a linear classifier and let \mathbf{q} be an adversarial perturbation for an instance \mathbf{x} where $f(\mathbf{x}) = c$ and $f(\mathbf{x} + \mathbf{q}) = c' \neq c$. There exists an $\alpha \in (0, 1)$ such that $\mathbf{w}_c \cdot (\mathbf{x} + \alpha\mathbf{q}) + b_c = \mathbf{w}_{c'} \cdot (\mathbf{x} + \alpha\mathbf{q}) + b_{c'}$.*

Proof. A little algebra shows

$$\alpha = \frac{(\mathbf{w}_{c'} - \mathbf{w}_c) \cdot \mathbf{x} + b_{c'} - b_c}{(\mathbf{w}_c - \mathbf{w}_{c'}) \cdot \mathbf{q}}. \quad (7)$$

It remains to be proven that $0 < \alpha < 1$. Note that, by assumption, $\mathbf{w}_c \cdot \mathbf{x} + b_c > \mathbf{w}_{c'} \cdot \mathbf{x} + b_{c'}$ and $\mathbf{w}_c \cdot (\mathbf{x} + \mathbf{q}) + b_c < \mathbf{w}_{c'} \cdot (\mathbf{x} + \mathbf{q}) + b_{c'}$. This implies $(\mathbf{w}_{c'} - \mathbf{w}_c) \cdot \mathbf{x} + b_{c'} - b_c < 0$ and $(\mathbf{w}_c - \mathbf{w}_{c'}) \cdot \mathbf{q} < 0$, so $\alpha > 0$. Furthermore, since $\mathbf{w}_c \cdot (\mathbf{x} + \mathbf{q}) + b_c < \mathbf{w}_{c'} \cdot (\mathbf{x} + \mathbf{q}) + b_{c'}$, we find $\alpha < 1$. \square

In other words, Lemma A.1 states that for any linear classifier f , if we want to adversarially perturb an input \mathbf{x} so its assigned label changes from c to c' , we must cross the boundary where f assigns equal probability to those two classes. Note that this condition is only necessary, not sufficient: it is perfectly possible to cross this boundary for the classes c and c' at a point where some other class c'' is still more likely than either c or c' , but we must cross this boundary nonetheless. Using Lemma A.1 we find

Lemma A.2. *Let f be a linear classifier and let \mathbf{x} be any input such that $f(\mathbf{x}) = c'$. Then for all classes $c \neq c'$,*

$$\Delta_{\text{adv}}(\mathbf{x}; f, c) \geq \frac{|(\mathbf{w}_{c'} - \mathbf{w}_c) \cdot \mathbf{x} + b_{c'} - b_c|}{\|\mathbf{w}_{c'} - \mathbf{w}_c\|}.$$

Proof. We can find the norm of the smallest perturbation \mathbf{r} such that $\mathbf{w}_c \cdot (\mathbf{x} + \mathbf{r}) + b_c = \mathbf{w}_{c'} \cdot (\mathbf{x} + \mathbf{r}) + b_{c'}$ using the following Lagrangian:

$$\mathcal{L} = \|\mathbf{r}\|^2 + \lambda((\mathbf{w}_{c'} - \mathbf{w}_c) \cdot (\mathbf{x} + \mathbf{r}) + b_{c'} - b_c)$$

The solution to this optimization problem is given by

$$\mathbf{r} = \frac{(\mathbf{w}_{c'} - \mathbf{w}_c) \cdot \mathbf{x} + b_{c'} - b_c}{\|\mathbf{w}_{c'} - \mathbf{w}_c\|^2} (\mathbf{w}_c - \mathbf{w}_{c'}).$$

Taking norms we get

$$\|\mathbf{r}\| = \frac{|(\mathbf{w}_{c'} - \mathbf{w}_c) \cdot \mathbf{x} + b_{c'} - b_c|}{\|\mathbf{w}_{c'} - \mathbf{w}_c\|}.$$

By Lemma A.1 and the construction of \mathbf{r} , any adversarial perturbation \mathbf{q} must satisfy $\|\mathbf{q}\| > \|\mathbf{r}\|$. In particular, any adversarial perturbation \mathbf{q} such that $f(\mathbf{x} + \mathbf{q}) = c$ satisfies $\|\mathbf{q}\| > \|\mathbf{r}\|$. Hence $\Delta_{\text{adv}}(\mathbf{x}; f, c) \geq \|\mathbf{r}\|$. \square

The following result is a trivial consequence of Lemma A.2:

Theorem A.3. *Let f be a linear classifier. Then for all inputs \mathbf{x} where $f(\mathbf{x}) = c$,*

$$\Delta_{\text{adv}}(\mathbf{x}; f) \geq \min_{c' \neq c} \frac{|(\mathbf{w}_{c'} - \mathbf{w}_c) \cdot \mathbf{x} + b_{c'} - b_c|}{\|\mathbf{w}_{c'} - \mathbf{w}_c\|}.$$

Note how Theorem A.3 confirms the intuition that the smallest adversarial perturbation to an instance \mathbf{x} is bounded from below by the orthogonal projection of \mathbf{x} onto the class closest to \mathbf{x} but distinct from $f(\mathbf{x})$, since this is exactly the quantity on the right-hand side of the inequality.

Theorem A.4. *The bound of Theorem A.3 is tight.*

Proof. Let f be a classifier for $C = d \geq 2$ classes where

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } x_1 > x_2, \dots, x_C \\ 2 & \text{if } x_2 > x_1, x_3, \dots, x_C \\ \vdots & \\ C & \text{if } x_C > x_1, \dots, x_{C-1} \end{cases}.$$

This classifier is linear since it can be characterized by linear functions of the form

$$f_i(\mathbf{x}) = \mathbf{e}_i \cdot \mathbf{x} = x_i.$$

Hence, it is subject to the bound of Theorem A.3, which in this case simplifies to

$$\Delta_{\text{adv}}(\mathbf{x}; f) \geq \min_{c' \neq c} |x_{c'} - x_c|.$$

It is easily seen, however, that this bound is exact for this particular classifier, i.e.

$$\Delta_{\text{adv}}(\mathbf{x}; f) = \min_{c' \neq c} |x_{c'} - x_c|.$$

The reason being that f classifies an input \mathbf{x} into the class corresponding to the index of the maximal component of \mathbf{x} . Thus, $f(\mathbf{x}) = c$ if and only if x_c is the maximal component of \mathbf{x} . To find the norm of the minimal perturbation \mathbf{r} such that $f(\mathbf{x} + \mathbf{r}) \neq c$, one simply takes $\mathbf{r} = (x_c - x_{c'})\mathbf{e}_{c'}$ where c' is the index of the second-highest component of \mathbf{x} . Clearly $f(\mathbf{x} + \mathbf{r}) = c' \neq c$ and $\|\mathbf{r}\| = |x_{c'} - x_c|$ is minimal. \square

A.2 Fully-connected layers

Applying Taylor's theorem to \mathbf{h}_L we obtain

$$\mathbf{h}_L(\mathbf{x} + \mathbf{r}) = \mathbf{h}_L(\mathbf{x}) + \mathbf{J}(\mathbf{x})\mathbf{r} + \boldsymbol{\varepsilon},$$

where

$$\|\boldsymbol{\varepsilon}\| \leq \frac{M}{2} \sqrt{n} \|\mathbf{r}\|^2.$$

Here, M is a real number bounding the absolute value of all second-order derivatives of \mathbf{h}_L from above. Thus $\mathbf{q} = \mathbf{J}(\mathbf{x})\mathbf{r} + \boldsymbol{\varepsilon}$ and

$$\|\mathbf{q}\| \leq \|\mathbf{J}(\mathbf{x})\| \|\mathbf{r}\| + \frac{M}{2} \sqrt{n} \|\mathbf{r}\|^2.$$

This is a quadratic inequality in $\|\mathbf{r}\|$ whose solution is given by

$$\|\mathbf{r}\| \geq \frac{\sqrt{\|\mathbf{J}(\mathbf{x})\|^2 + 2M\sqrt{n}\|\mathbf{q}\|} - \|\mathbf{J}(\mathbf{x})\|}{M\sqrt{n}}.$$

By Theorem A.3 we know

$$\|\mathbf{q}\| \geq \min_{c' \neq c} \frac{|(\mathbf{w}_c - \mathbf{w}_{c'}) \cdot \mathbf{h}(\mathbf{x})|}{\|\mathbf{w}_c - \mathbf{w}_{c'}\|}.$$

Hence the theorem follows.

One might rightly ask how realistic this result actually is. After all, we needed to assume that the function \mathbf{h}_L was twice differentiable and had bounded second-order derivatives. In this section, we will try to show that these assumptions are quite realistic by deriving them from other realistic assumptions on the activation functions. Specifically, we have the following

Theorem A.5. Consider fully-connected layers $\mathbf{h}_1, \dots, \mathbf{h}_L$ whose activation functions g_i are twice differentiable on \mathbb{R} . Assume also that for each g_i there exist $N_i > 0$ and $M_i > 0$ such that $|g'_i(x)| \leq N_i$ and $|g''_i(x)| \leq M_i$ for all x . Then $\mathbf{h}_L(\mathbf{x}) = [h_L^{(1)}(\mathbf{x}), \dots, h_L^{(n)}(\mathbf{x})]^T$ satisfies the following properties:

1. each $h_L^{(i)}$ is twice differentiable;
2. $|\partial^\alpha h_L^{(i)}(\mathbf{x})| \leq M$ for all i , $|\alpha| = 2$ and some $M > 0$;
3. $|\partial^\alpha h_L^{(i)}(\mathbf{x})| \leq N$ for all i , $|\alpha| = 1$ and some $N > 0$.

Proof. The proof proceeds by induction on L . The case where $L = 0$ is trivial, so suppose

1. each $h_L^{(i)}$ is twice differentiable;
2. $|\partial^\alpha h_L^{(i)}(\mathbf{x})| \leq M'$ for all i , $|\alpha| = 2$ and some $M' > 0$;
3. $|\partial^\alpha h_L^{(i)}(\mathbf{x})| \leq N'$ for all i , $|\alpha| = 1$ and some $N' > 0$.

We need to show that $\mathbf{h}_{L+1} : \mathbb{R}^d \rightarrow \mathbb{R}^n : \mathbf{x} \mapsto [h_{L+1}^{(1)}(\mathbf{x}), \dots, h_{L+1}^{(n)}(\mathbf{x})]^T$ then satisfies the following properties:

1. $h_{L+1}^{(i)}$ is twice differentiable for all i ;
2. there exists an $M > 0$ such that $|\partial^\alpha h_{L+1}^{(i)}(\mathbf{x})| \leq M$ for all \mathbf{x} , $|\alpha| = 2$ and i ;
3. there exists an $N > 0$ such that $|\partial^\alpha h_{L+1}^{(i)}(\mathbf{x})| \leq N$ for all \mathbf{x} , $|\alpha| = 1$ and i .

Since $\mathbf{h}_{L+1}(\mathbf{x}) = g_{L+1}(\mathbf{V}_{L+1}\mathbf{h}_L(\mathbf{x}) + \mathbf{b}_{L+1})$ we find

$$h_{L+1}^{(i)}(\mathbf{x}) = g_{L+1} \left(\sum_j v_{L+1,i,j} h_L^{(j)}(\mathbf{x}) + b_{L+1,i} \right).$$

Clearly, since g_{L+1} is twice differentiable by assumption and $h_L^{(j)}$ is twice differentiable for all j by the induction hypothesis, $h_{L+1}^{(i)}$ is twice differentiable for all i . This shows Item 1. To show Item 2, we distinguish two cases. First, let

$$\alpha_j = \begin{cases} 2 & j = k \\ 0 & \text{otherwise} \end{cases}$$

for some $k \in \{1, \dots, d\}$. Then

$$\begin{aligned} \partial^\alpha h_{L+1}^{(i)} &= \frac{\partial^2 h_{L+1}^{(i)}}{\partial x_k^2} = \frac{\partial}{\partial x_k} \left(g'_{L+1} \left(\sum_j v_{L+1,i,j} h_L^{(j)}(\mathbf{x}) + b_{L+1,i} \right) \sum_j v_{L+1,i,j} \frac{\partial}{\partial x_k} h_L^{(j)}(\mathbf{x}) \right) \\ &= g''_{L+1} \left(\sum_j v_{L+1,i,j} h_L^{(j)}(\mathbf{x}) + b_{L+1,i} \right) \left(\sum_j v_{L+1,i,j} \frac{\partial}{\partial x_k} h_L^{(j)}(\mathbf{x}) \right)^2 + \\ &\quad g'_{L+1} \left(\sum_j v_{L+1,i,j} h_L^{(j)}(\mathbf{x}) + b_{L+1,i} \right) \sum_j v_{L+1,i,j} \frac{\partial^2}{\partial x_k^2} h_L^{(j)}(\mathbf{x}). \end{aligned}$$

Taking absolute values and applying the induction hypothesis, this yields

$$|\partial^\alpha h_{L+1}^{(i)}| \leq M_{L+1} \left(N' \sum_j |v_{L+1,i,j}| \right)^2 + N_{L+1} M' \sum_j |v_{L+1,i,j}|.$$

Hence we may choose

$$M = M_{L+1} (N' s)^2 + N_{L+1} M' s,$$

where

$$s = \max_i \sum_j |v_{L+1,i,j}|.$$

For the second case, let

$$\alpha_j = \begin{cases} 1 & j \in \{k_1, k_2\} \\ 0 & \text{otherwise} \end{cases}$$

for $k_1, k_2 \in \{1, \dots, d\}$ and $k_1 < k_2$. Of course, this case only applies when $d > 1$. We find

$$\begin{aligned} \partial^\alpha h_{L+1}^{(i)} &= \frac{\partial^2 h_{L+1}^{(i)}}{\partial x_{k_1} \partial x_{k_2}} = \frac{\partial}{\partial x_{k_2}} \left(g'_{L+1} \left(\sum_j v_{L+1,i,j} h_L^{(j)}(\mathbf{x}) + b_{L+1,i} \right) \sum_j v_{L+1,i,j} \frac{\partial}{\partial x_{k_1}} h_L^{(j)}(\mathbf{x}) \right) \\ &= g''_{L+1} \left(\sum_j v_{L+1,i,j} h_L^{(j)}(\mathbf{x}) + b_{L+1,i} \right) \sum_j v_{L+1,i,j} \frac{\partial}{\partial x_{k_2}} h_L^{(j)}(\mathbf{x}) \sum_j v_{L+1,i,j} \frac{\partial}{\partial x_{k_1}} h_L^{(j)}(\mathbf{x}) \\ &\quad + g'_{L+1} \left(\sum_j v_{L+1,i,j} h_L^{(j)}(\mathbf{x}) + b_{L+1,i} \right) \sum_j v_{L+1,i,j} \frac{\partial^2}{\partial x_{k_1} \partial x_{k_2}} h_L^{(j)}(\mathbf{x}) \end{aligned}$$

Again, taking absolute values and applying the induction hypothesis:

$$|\partial^\alpha h_{L+1}^{(i)}(\mathbf{x})| \leq M_{L+1} \left(N' \sum_j |v_{L+1,i,j}| \right)^2 + N_{L+1} M' \sum_j |v_{L+1,i,j}|$$

The result is identical to the first case.

Finally, to show Item 3, we let

$$\alpha_j = \begin{cases} 1 & j = k \\ 0 & \text{otherwise} \end{cases}$$

for some $k \in \{1, \dots, d\}$. Then

$$\partial^\alpha h_{L+1}^{(i)} = \frac{\partial h_{L+1}^{(i)}}{\partial x_k} = g'_{L+1} \left(\sum_j v_{L+1,i,j} h_L^{(j)}(\mathbf{x}) + b_{L+1,i} \right) \sum_j v_{L+1,i,j} \frac{\partial}{\partial x_k} h_L^{(j)}(\mathbf{x}).$$

Again taking absolute values and applying the induction hypothesis, we find

$$|\partial^\alpha h_{L+1}^{(i)}| \leq N_{L+1} N' \sum_j |v_{L+1,i,j}|.$$

Hence we may choose

$$N = (N_{L+1} N') \max_i \sum_j |v_{L+1,i,j}|.$$

This completes the proof. \square

By Theorem A.5, in order for Theorem 4.1 to hold it is sufficient that the activation functions of the MLP in question be twice differentiable and have bounded first and second derivatives. These assumptions are not unrealistic: they are satisfied by the logistic sigmoid function, for example. The logistic sigmoid is in fact just a scaled and shifted version of the hyperbolic tangent:

$$\text{sigm}(x) = \frac{1}{2} \tanh\left(\frac{x}{2}\right) + \frac{1}{2}. \quad (8)$$

Since \tanh is twice differentiable, so is sigm . Moreover, the first and second derivatives of \tanh are bounded by 1, so the first and second derivatives of sigm are bounded as well. In fact, $|\tanh(x)| \leq 1$ for all x .

The ReLU activation function presents some problems, however, as it is not differentiable at zero. Gradient-based optimization requires all activation functions be differentiable, though, so in practice either a smooth approximation to ReLU is used which is differentiable, such as the softplus function $\ln(1+\exp(x))$, or the value of the derivative is simply set to zero at the origin (Glorot et al. [2011]). In both cases it can be seen that ReLU (as it is used in practice) also satisfies the necessary assumptions for Theorem 4.1 to hold.

Note also how the proof of Theorem A.5 yields an efficient algorithm for approximating the M parameter used in the lower bound. Algorithm A.1 shows how this can be done in $\mathcal{O}(n)$ time where n is the number of parameters of the neural network.

<p>Algorithm A.1: Computation of M</p> <p>Data: MLP f with L hidden layers, activation functions g_i satisfying $g'_i(x) \leq A_i$ and $g''_i(x) \leq B_i$ for all x.</p> <p>Result: a value M satisfying $\partial^\alpha h_L^{(i)}(\mathbf{x}) \leq M$ for all \mathbf{x}, $\alpha = 2$ and i</p> <p>begin</p> <p style="padding-left: 1em;">$M_0 \leftarrow 0$</p> <p style="padding-left: 1em;">$N_0 \leftarrow 1$</p> <p style="padding-left: 1em;">for i <i>from</i> 1 <i>to</i> L do</p> <p style="padding-left: 2em;">$s \leftarrow \max_j \sum_k v_{i,j,k}$</p> <p style="padding-left: 2em;">$M_i \leftarrow B_i N_{i-1}^2 s^2 + A_i M_{i-1} s$</p> <p style="padding-left: 2em;">$N_i \leftarrow A_i N_{i-1} s$</p> <p style="padding-left: 1em;">end</p> <p style="padding-left: 1em;">return M_L</p> <p>end</p>

A.3 Convolutional layers

Using Lemma B.1 we find¹

$$\begin{aligned} \|\text{ReLU}(\mathbf{W} \star (\mathbf{X} + \mathbf{R}) + \mathbf{b})\|_F &= \|\text{ReLU}(\mathbf{W} \star \mathbf{X} + \mathbf{W} \star \mathbf{R} + \mathbf{b})\|_F \\ &\leq \|\text{ReLU}(\mathbf{W} \star \mathbf{X} + \mathbf{b}) + \text{ReLU}(\mathbf{W} \star \mathbf{R})\|_F \\ &\leq \|\text{ReLU}(\mathbf{W} \star \mathbf{X} + \mathbf{b})\|_F + \|\text{ReLU}(\mathbf{W} \star \mathbf{R})\|_F. \end{aligned}$$

This yields

$$\|\text{ReLU}(\mathbf{W} \star \mathbf{R})\|_F \geq \kappa.$$

A necessary condition for this equality to hold is (Lemma B.1)

$$\|\mathbf{W} \star \mathbf{R}\|_F = \kappa. \tag{9}$$

Thanks to Lemma B.3, we may rewrite Equation (9) as

$$\|\mathbf{R}\|_F \geq \frac{\kappa}{\|\mathbf{W}\|_F}.$$

A.4 Pooling layers

Assuming any adversarial perturbation \mathbf{Q} to the input of the next layer needs to satisfy $\|\mathbf{Q}\|_F \geq \kappa$, Assumption 4.3 implies we have to solve the following equation:

$$\|\mathbf{Z}(\mathbf{R})\|_F \geq \kappa. \tag{10}$$

¹Note that even though $\mathbf{W} \star (\mathbf{X} + \mathbf{R}) + \mathbf{b}$ does not represent a real convolution in this context, the linearity property of Lemma B.1 still applies.

A necessary condition for Equation (10) to hold is to have

$$|z_{ijk}(\mathbf{R})| \geq \frac{\kappa}{t} \quad (11)$$

for at least one element $z_{ijk}(\mathbf{R})$. How this can be done depends on the precise nature of the pooling operation.

A.4.1 MAX-pooling

MAX-pooling reduces the dimensionality of the input by taking the maximum of all $q \times q$ regions within the receptive field:

$$z_{ijk}(\mathbf{X}) = \max\{x_{inm} \mid (n, m) \in I(j, k)\}.$$

Proof that MAX-pooling satisfies Assumption 4.3 is given in Lemma B.4. In order to satisfy Equation (11), it is clearly necessary to set at least one component of \mathbf{R} equal to or greater than κ/t in absolute value. This yields

$$\|\mathbf{R}\|_F \geq \frac{\kappa}{t}. \quad (12)$$

A.4.2 L_p pooling

An L_p pooling layer produces as output the L_p norm of its input:

$$z_{ijk}(\mathbf{X}) = \left(\sum_{(n,m) \in I} |x_{inm}|^p \right)^{\frac{1}{p}}.$$

We will write \mathbf{v}_{ijk} for the vector whose L_p norm is taken in the computation of z_{ijk} . Proof that L_p -pooling satisfies Assumption 4.3 is given in Lemma B.5. In order to satisfy Equation (11), we must have

$$\|\mathbf{v}_{ijk}(\mathbf{R})\|_p \geq \frac{\kappa}{t} \quad (13)$$

for some i, j, k . For Equation (13) to be satisfied, there must be at least one element r_{lmn} in some receptive field of \mathbf{R} such that

$$|r_{lmn}| \geq \frac{\kappa}{tq^{2/p}}.$$

Since there will be at least one receptive field, at least one element of \mathbf{R} must satisfy this requirement and hence

$$\|\mathbf{R}\|_F \geq \frac{\kappa}{tq^{2/p}}. \quad (14)$$

Note the nice property that as $p \rightarrow \infty$ we find

$$\|\mathbf{R}\|_F \geq \frac{\kappa}{t},$$

which is the bound for MAX-pooling.

A.4.3 Average pooling

An average pooling layer takes the average of all its inputs:

$$z_{ijk}(\mathbf{X}) = \frac{1}{q^2} \sum_{(n,m) \in I} x_{inm}.$$

Proof that average pooling satisfies Assumption 4.3 is given in Lemma B.6. Note that, contrary to all the other pooling operations studied here, Assumption 4.3 holds with equality in the case of average pooling. To satisfy Equation (11), it is necessary that at least one element of \mathbf{R} be greater than or equal to κ/t in absolute value. We thus find

$$\|\mathbf{R}\|_F \geq \frac{\kappa}{t}. \quad (15)$$

B Auxiliary results

Lemma B.1.

1. Let $a, b \in \mathbb{R}$, then

$$\text{ReLU}(a + b) \leq \text{ReLU}(a) + \text{ReLU}(b).$$

2. Let $\mathbf{W}, \mathbf{X}, \mathbf{R}$ be tensors and $s \in \mathbb{N}$, then

$$\mathbf{W} \star_s (\mathbf{X} + \mathbf{R}) = \mathbf{W} \star_s \mathbf{X} + \mathbf{W} \star_s \mathbf{R}.$$

3. Let \mathbf{X} be any real-valued tensor, then

$$\|\text{ReLU}(\mathbf{X})\|_F \leq \|\mathbf{X}\|_F.$$

Proof.

1. We distinguish four cases:

- $a, b > 0$:

$$\text{ReLU}(a + b) = a + b = \text{ReLU}(a) + \text{ReLU}(b).$$

- $a > 0$ and $b \leq 0$:

$$\text{ReLU}(a + b) \leq a + b \leq a = \text{ReLU}(a) + \text{ReLU}(b).$$

- $a \leq 0$ and $b > 0$:

$$\text{ReLU}(a + b) \leq a + b \leq b = \text{ReLU}(a) + \text{ReLU}(b).$$

- $a, b < 0$:

$$\text{ReLU}(a + b) = 0 = \text{ReLU}(a) + \text{ReLU}(b).$$

2. Suppose $\mathbf{W} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ and $\mathbf{X}, \mathbf{R} \in \mathbb{R}^{m_1 \times \dots \times m_d}$, then

$$\begin{aligned} (\mathbf{W} \star_s (\mathbf{X} + \mathbf{R}))_{i_1 \dots i_d} &= \sum_{j_1, \dots, j_d} w_{j_1 \dots j_d} (x_{i_1 + s(j_1 - 1), \dots, i_d + s(j_d - 1)} + r_{i_1 + s(j_1 - 1), \dots, i_d + s(j_d - 1)}) \\ &= \sum_{j_1, \dots, j_d} w_{j_1 \dots j_d} x_{i_1 + s(j_1 - 1), \dots, i_d + s(j_d - 1)} + \sum_{j_1, \dots, j_d} w_{j_1 \dots j_d} r_{i_1 + s(j_1 - 1), \dots, i_d + s(j_d - 1)} \\ &= (\mathbf{W} \star_s \mathbf{X})_{i_1 \dots i_d} + (\mathbf{W} \star_s \mathbf{R})_{i_1 \dots i_d}. \end{aligned}$$

3. Let $\mathbf{X} \in \mathbb{R}^{n_1 \times \dots \times n_d}$. We compute:

$$\|\text{ReLU}(\mathbf{X})\|_F^2 = \sum_{i_1, \dots, i_d} \max\{0, x_{i_1 \dots i_d}\}^2 \leq \sum_{i_1, \dots, i_d} x_{i_1 \dots i_d}^2 = \|\mathbf{X}\|_F^2.$$

□

Lemma B.2. Let a_i, b_i be non-negative real numbers for $i = 1, \dots, n$. Then

$$\sum_i a_i b_i \leq \left(\sum_i a_i \right) \left(\sum_i b_i \right).$$

Proof. This is easy to see by direct computation:

$$\left(\sum_i a_i \right) \left(\sum_i b_i \right) = \sum_i a_i \left(\sum_j b_j \right).$$

Since all terms are non-negative, we have

$$b_i \leq \sum_j b_j$$

for all i . Hence the result follows. □

Lemma B.3. Let \mathbf{W} and \mathbf{R} be as above, then

$$\|\mathbf{W} \star \mathbf{R}\|_F \leq \|\mathbf{W}\|_F \|\mathbf{R}\|_F.$$

Proof. We compute:

$$\begin{aligned} \|\mathbf{W} \star \mathbf{R}\|_F^2 &= \sum_i \|\mathbf{W}_i \star \mathbf{R}\|_F^2 = \sum_{i,j,k} (\mathbf{W}_i \star \mathbf{R})_{jk}^2 \\ &= \sum_{i,j,k,l,m,n} w_{ilmn}^2 r_{l,m+s(j-1),n+s(k-1)}^2. \end{aligned}$$

Using Lemma B.2:

$$\begin{aligned} \sum w_{ilmn}^2 r_{l,m+s(j-1),n+s(k-1)}^2 &\leq \left(\sum w_{ilmn}^2 \right) \left(\sum r_{l,m+s(j-1),n+s(k-1)}^2 \right) \\ &\leq \left(\sum w_{ilmn}^2 \right) \left(\sum r_{lmn}^2 \right) \\ &= \|\mathbf{W}\|_F^2 \|\mathbf{R}\|_F^2. \end{aligned}$$

We can thus conclude

$$\|\mathbf{W} \star \mathbf{R}\|_F \leq \|\mathbf{W}\|_F \|\mathbf{R}\|_F. \quad \square$$

Lemma B.4. MAX-pooling satisfies Assumption 4.3.

Proof. We compute:

$$\begin{aligned} z_{ijk}(\mathbf{X} + \mathbf{R}) &= \max\{x_{inm} + r_{inm} \mid (n, m) \in I(j, k)\} \\ &\leq \max\{x_{inm} \mid (n, m) \in I(j, k)\} + \max\{r_{inm} \mid (n, m) \in I(j, k)\} \\ &= z_{ijk}(\mathbf{X}) + z_{ijk}(\mathbf{R}). \end{aligned} \quad \square$$

Lemma B.5. L_p pooling satisfies Assumption 4.3.

Proof. Define $\mathbf{v}_{ijk}(\mathbf{X})$ to be the vector whose L_p norm is taken in the computation of $z_{ijk}(\mathbf{X})$. We find

$$\begin{aligned} z_{ijk}(\mathbf{X} + \mathbf{R}) &= \|\mathbf{v}_{ijk}(\mathbf{X} + \mathbf{R})\|_p = \|\mathbf{v}_{ijk}(\mathbf{X}) + \mathbf{v}_{ijk}(\mathbf{R})\|_p \\ &\leq \|\mathbf{v}_{ijk}(\mathbf{X})\|_p + \|\mathbf{v}_{ijk}(\mathbf{R})\|_p = z_{ijk}(\mathbf{X}) + z_{ijk}(\mathbf{R}). \end{aligned} \quad \square$$

Lemma B.6. Average pooling satisfies Assumption 4.3.

Proof. We have

$$\begin{aligned} z_{ijk}(\mathbf{X} + \mathbf{R}) &= \frac{1}{q^2} \sum_{(n,m) \in I} (x_{inm} + r_{inm}) \\ &= \frac{1}{q^2} \sum_{(n,m) \in I} x_{inm} + \frac{1}{q^2} \sum_{(n,m) \in I} r_{inm} \\ &= z_{ijk}(\mathbf{X}) + z_{ijk}(\mathbf{R}). \end{aligned} \quad \square$$