Homework 4

1. Architectural choice v.s Vanishing/Exploding Gradients.

As per Piazza @761. this question considers a N layer MLP with a single unit in all the hidden layers and the weight matrices are set to 1.

### 1.1.1 Effect of Activation - Sigmoid [1pt]

$$\tau(z) = \frac{1}{1+e^{-z}}.$$

$$\tau'(z) = \frac{-1}{(1+e^{-z})^2}(-e^{-z})$$

$$= \frac{e^{-z}}{(1+e^{-z})^2}$$

$$\max(\tau'(z)) = 0.25 = \frac{1}{4}.$$

$$\min(\tau'(z)) = 0.$$

$$\left|\frac{\partial f(x)}{\partial x}\right| = \left|\frac{\partial h^n}{\partial x}\right| = \left|\frac{\partial h^n}{\partial h^{n-1}}\frac{\partial h^{n-1}}{\partial x}\right| = \left|\tau'(h^{n-1})\frac{\partial h^{n-1}}{\partial x}\right|$$

$$0 \le \left|\frac{\partial f(x)}{\partial x}\right| \le \frac{1}{4}\left|\frac{\partial h^{n-1}}{\partial x}\right|$$

$$\le \frac{1}{4^2}\left|\frac{\partial h^{n-2}}{\partial x}\right|$$

$$\le \cdots$$

$$\le \frac{1}{4}^n$$

therefore, $0 \le \left|\frac{\partial f(x)}{\partial x}\right| \le \frac{1}{4}^n$.

As depth n approaches $\infty$, the gradients will necessarily vanish but not explode.

### 1.1.2 Effect of Activation - Tanh [1pt]

$$\tau(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\tau'(z) = 1 - \tau(z)^2$$

$$= 1 - \frac{e^{2z} - 2 + e^{-2z}}{e^{2z} + 2 + e^{-2z}}$$

$$= \frac{e^{2z} + 2 + e^{-2z} - e^{2z} + 2 - e^{-2z}}{e^{2z} + 2 + e^{-2z}}$$

$$= \frac{4}{e^{2z} + 2 + e^{-2z}}$$

$$\max(\tau'(z)) = 1$$

$$\min(\tau'(z)) = 0$$

Similarly, $0 \le \left|\frac{\partial f(x)}{\partial x}\right| \le 1\left|\frac{\partial h^{n-1}}{\partial x}\right|$

$$\le \cdots$$

$$\le 1^n = 1$$

As depth n approaches $\infty$, the gradients don't necessarily explode or vanish since it can remain unchanged through backprop.

### 1.2.1 Gradient through RNN [1pt]

$$T_{max}\left(\frac{\partial h_n}{\partial x_1}\right) \le T_{max}\left(\frac{\partial h_n}{\partial h_{n-1}}\right)T_{max}\left(\frac{\partial h_{n-1}}{\partial x_1}\right).$$

$$\le T_{max}\left(\frac{\partial h_n}{\partial h_{n-1}}\right)T_{max}\left(\frac{\partial h_{n-1}}{\partial h_{n-2}}\right)T_{max}\left(\frac{\partial h_{n-2}}{\partial x_1}\right)$$

$$\le \cdots$$

$$\le T_{max}\left(\frac{\partial h_n}{\partial h_{n-1}}\right)\cdots T_{max}\left(\frac{\partial h_2}{\partial x_1}\right)$$

$$\frac{\partial h_{t+1}}{\partial x_t} = \frac{\partial h_{t+1}}{\partial W x_t}\frac{\partial W x_t}{\partial x_t} = diag\left(\frac{4}{e^{2z_t} + 2 + e^{-2z_t}}\right)W.$$

where $z_t = W x_t$

$$T_{max}\left(\frac{\partial h_{t+1}}{\partial x_t}\right) \le T_{max}\left(diag\left(\frac{4}{e^{2z_t} + 2 + e^{-2z_t}}\right)T_{max}(W)\right)$$

$$\le T_{max}\left(\frac{\partial x_{t+1}}{\partial x_t}\right)T_{max}(W) = \frac{1}{2}$$

Substitute back;

$$0 \le T_{max}\left(\frac{\partial h_n}{\partial x_1}\right) \le \frac{1}{2}^n$$

### 1.3.2 Batch Normalization and ResNet. [1pt]

the one on the **left** is easier to learn because its gradient doesn't vanish,

Looking at one such block on the left

$$\frac{\partial h_k}{\partial h_{k-1}} = (1 + grad through block).$$

there will always be a 1 added to the gradient, passing through unchanged, This is not the case for the architecture on the right.

### 1.2.3 Benefits of Residual Connections [1pt]

Result from 1.2.2 states that:

$$T_{min}\left(\frac{\partial z_{t+1}}{\partial z_t}\right) \ge 1 - \delta_{small}.$$

$$T_{max}\left(\frac{\partial z_{t+1}}{\partial z_t}\right) \ge \delta_{big} - 1., \quad \delta_{big} >> 2$$

So, $T_{max}\left(\frac{\partial z_{t+1}}{\partial z_t}\right) >> 1$.

Similarly to 1.2.1,

$$T_{max}\left(\frac{\partial z_n}{\partial z_1}\right) \le T_{max}\left(\frac{\partial z_n}{\partial z_{n-1}}\right)\cdots T_{max}\left(\frac{\partial z_2}{\partial z_1}\right).$$

will not vanish, since products of terms >> 1 will never be 0.

However, it doesn't solve the exploding gradient problem, since the product can still approach $\infty$.

2. Autoregressive Models.
2.2 PixelCNN
### 2.2.1 Connections [1pt].
$O(WHdk^2)$.

### 2.2.2 Parallelism [1pt].
$O(d)$

2.3 Multidimensional RNN.
### 2.3.1 Connections [1pt]
$O(WHdk)$.

Fundamentally in terms of computational complexity MDRNN is better since it's $O(WHdk)$ # of connections whereas PixelCNN is $O(WHdk^2)$

### 2.3.3 Discussion [1pt].
PixelCNN is better in terms of parallelization, As discussed before it's sequential operation is $O(d)$, whereas MDRNN neurons' computation are not independent, so they cannot be computed in parallel.