### 1.1.1. Sigmoid

$x \to h_1 \to \cdots \to h_n \to y.$

$h_1 = \sigma(w_1 x + b_1)$

$h_2 = \sigma(w_2 h_1 + b_2)$

$\cdots$

$y = \sigma(w_n h_{n-1} + b_n)$

$w_1 = w_2 = \cdots = w_n = 1$

if $n = 1$. $\quad y = \sigma(w_1 x + b_1) = \sigma(x + b_1)$

$\dfrac{dy}{dx} = \sigma(x + b_1)(1 - \sigma(x + b_1))$

$= h_1 \cdot (1 - h_1)$

$\because h_1$ is a sigmoid func, $0 < h_1 < 1$. $h_1'$ is max when $x = \frac{1}{2}$

$\therefore 0 \le \left|\dfrac{dy}{dx}\right| \le \frac{1}{4} \Rightarrow |h_1'| \le \frac{1}{4}$

for $n = i$. $\quad y = \sigma(h_{i-1} + b_{i-1}) = h_i$.

$\dfrac{dy}{dx} = \sigma(h_{i-1} + b_{i-1})(1 - \sigma(h_{i-1} + b_{i-1})) \cdot h_{i-1}'$

$0 \le \left|\dfrac{dy}{dx}\right| = \left|\dfrac{dh_i}{dx}\right| \le \frac{1}{4} h_{i-1}'$

$\Rightarrow$ By induction $\Rightarrow$ $0 \le \left|\dfrac{dy}{dx}\right| = \left|\dfrac{dh_n}{dx}\right| \le \left(\frac{1}{4}\right)^n$

For $\dfrac{dh_1}{dx}$. we have $0 \le \left|\dfrac{dh_1}{dx}\right| \le \frac{1}{4}$

$\Rightarrow \left|\dfrac{dh_n}{dh_1}\right| = \left|\dfrac{dh_n}{dx} \cdot \dfrac{dx}{dh_1}\right| \Rightarrow 0 \le \dfrac{dh_n}{dh_1} \le \left(\frac{1}{4}\right)^{n-1}$ as $n \to \infty$. $\left(\frac{1}{4}\right)^{n-1} \to 0$

$\lim\limits_{n \to \infty} 0 \le \lim\limits_{n \to \infty}\left|\dfrac{dh_n}{dh_1}\right| \le \lim\limits_{n \to \infty}\left(\frac{1}{4}\right)^{n-1} \Rightarrow 0 \le \lim\limits_{n \to \infty}\left|\dfrac{dh_n}{dh_1}\right| \le 0$

$\Rightarrow \lim\limits_{n \to \infty}\left|\dfrac{dh_n}{dh_1}\right| = 0 \Rightarrow$ The gradient vanishes for Sigmoid activation

### 1.1.2. Tanh.

$x \to h_1 \to h_2 \cdots \to h_n \to y.$

$h_1 = \tanh(x + b_1)$

$\cdots$

$y = \tanh(h_{n-1} + b_{n-1})$

$n = 1$

$\dfrac{dh_1}{dx} = 1 - \tanh^2(x).$

$-1 < \tanh(x) < 1$

$\Rightarrow 0 \le \tanh(x) < 1$

$0 < \left|\dfrac{dh_1}{dx}\right| \le 1$

$n = i$ $\quad \dfrac{dh_i}{dx} = h_{i-1}'(1 - \tanh^2(h_{i-1} x + b_{i-1}))$

$0 < \left|\dfrac{dh_i}{dx}\right| \le h_{i-1}'$

By induction

$\Rightarrow 0 < \left|\dfrac{dy}{dx}\right| \le h_1' \Rightarrow 0 < \left|\dfrac{dy}{dx}\right| \le 1$

$0 < \left|\dfrac{dh_n}{dh_1}\right| = \left|\dfrac{dh_n}{dx} \cdot \dfrac{dx}{dh_{n-1}}\right| \le 1$

$\Rightarrow$ The gradient does not vanish or explode for Tanh function.

## 1.2.1.

$$\sigma_{max}\left(\frac{\partial x_n}{\partial x_1}\right) = \sigma_{max}\left(\frac{\partial x_n}{\partial x_{n-1}} \cdots \frac{\partial x_2}{\partial x_1}\right)$$

$$\leq \sigma_{max}\left(\frac{\partial x_n}{\partial x_{n-1}}\right) \cdots \sigma_{max}\left(\frac{\partial x_2}{\partial x_1}\right)$$

For any $i$ in $n \cdots 2$.

$$\frac{\partial x_i}{\partial x_{i-1}} = \frac{\partial \tanh(W x_{i-1})}{\partial x_{i-1}} = W(1 - \tanh^2(W x_{i-1}))$$

$$\sigma_{max}\left(\frac{\partial x_i}{\partial x_{i-1}}\right) = \sigma_{max}(W(1 - \tanh^2(W x_{i-1}))) \leq \sigma_{max}(W)\,\sigma_{max}(1 - \tanh^2(W x_{i-1}))$$

$$\because \sigma_{max}(W) = \tfrac{1}{2}. \quad \sigma_{max}(0) < \sigma_{max}(1 - \tanh^2(W x_{i-1})) \leq \sigma_{max}(1)$$

$$\Rightarrow \quad 0 < \sigma_{max}(1 - \tanh^2(W x_{i-1})) \leq 1$$

$$\therefore \quad 0 < \sigma_{max}\left(\frac{\partial x_i}{\partial x_{i-1}}\right) \leq \tfrac{1}{2}$$

$$\Rightarrow \quad 0 \leq \sigma_{max}\left(\frac{\partial x_n}{\partial x_1}\right) \leq \left(\tfrac{1}{2}\right)^n$$

## 1.2.3.

We know that (from 1.2.1)

$$\begin{cases} \sigma_{min}\left(\frac{z_{t+1}}{z_t}\right) \geq 1 - \sigma_{small} \\ \sigma_{max}\left(\frac{z_{t+1}}{z_t}\right) \geq \sigma_{big} - 1 \end{cases} \quad \text{and} \quad \begin{cases} \sigma_{small} \ll 1 \\ \sigma_{big} \gg 2. \end{cases}$$

$\Rightarrow$ Similarly to 1.2.1 $\qquad z_{t+1} = z_t + f_t \quad \frac{dz_{t+1}}{dz_t} = 1 + \frac{df_t}{dz_t}$

$$\sigma_{min}\left(\frac{\partial z_n}{\partial z_1}\right) = \sigma_{min}\left(\frac{\partial z_n}{\partial z_{n-1}} \cdots \frac{\partial z_2}{\partial z_1}\right) \geq \sigma_{min}\left(\frac{\partial z_n}{\partial z_{n-1}}\right) \cdots \sigma_{min}\left(\frac{\partial z_2}{\partial z_1}\right)$$

$$\geq (1 - \sigma_{small})^{n-1}$$

Since $\sigma_{small} \ll 1$

with $n < \infty$

$\Rightarrow \sigma_{min}\left(\frac{\partial z_n}{\partial z_1}\right) \geq 1$, is bounded, close to 1.   not explode/vanish.

$$\sigma_{max}\left(\frac{\partial z_n}{\partial z_1}\right) \leq \sigma_{max}\left(\frac{\partial z_n}{\partial z_{n-1}}\right) \cdots \sigma_{max}\left(\frac{\partial z_2}{\partial z_1}\right) \geq (\sigma_{big} - 1)^{n-1}$$

Since $\sigma_{big} \gg 2$

$\sigma_{big} - 1 \gg 1$

$(\sigma_{big} - 1)^{n-1} \rightsquigarrow \infty \quad \Rightarrow \sigma_{max}\left(\frac{\partial z_n}{\partial z_1}\right) \to \infty$, not bounded, explodes.

## 1.3.2.

The first model on the left is favoured

Because, for it $i$th layer and $i+1$ th:

$$\bar{x}^i = \bar{x}^{i+1}\left(1 + \frac{\partial F}{\partial x}\right). \text{ this model is more stable with a "1" in it,}$$

$$\text{to anti-vanishing to } 0.$$

2.2.1

$O(WHdk^2)$

2.2.2.

$O(d)$



2.3.1    $O(WHdk^2)$

2.3.3.

|  | Pixel CNN | MDRNN |
|---|---|---|
| Computational Memory. | weight: $4k^2d$   $\Rightarrow O(k^2d)$ | Linear in number of data points and the number of network weights. weight: $(2+2d)dk^2WH$   $\Rightarrow O(WHdk^2)$ |
| Computational Complexity. | $\approx$ connections.   $O(WHdk^2)$ | $O(WHdk^2)$ |
| Parallelism. | $O(d)$. parallel in each layer | Serial in each layer. |
| Size of Context window. | all data unmasked. | adj 2 hidden units. |

Conclusion: With similar computational task, pixel CNN is more suitable in parallelism, with less computation memory needed and larger size of context window; while MDRNN may capture more sequential relationships between data.