**Exercise 3**

Show that a multi-layer neural network with linear activation function $s(x) = x$ is equivalent to a single layer linear network. Assume that in each layer the inputs follow a Normal distribution with mean zero and small variance, *i.e.* $\sigma \ll 1$. For which of the activation functions $s(x) = 1/(1 + \exp(-x))$, $s(x) = \tanh(x)$, $s(x) = \text{relu}(x)$ and $s(x) = \text{selu}(x)$ is a deep network equivalent to a linear network for the given distribution ? The selu function is given by:

$$\text{selu}(x) = \begin{cases} \lambda x \text{ if } x > 0 \\ \alpha(\exp(x) - 1) \text{ otherwise} \end{cases}, \tag{4}$$

where $\lambda \approx 1.0507$ and $\alpha \approx 1.75814$. (*Hint:* consider the case $\sigma \to 0$ using a Taylor expansion around 0.)

- Deep Linear = Single linear

  MLP (with L layers) has about: $y = W^L \cdot W^{L-1} \cdot \ldots W^2 \cdot W^1 \cdot x$

  We can just taunle $W^L \cdot W^{L-1} \ldots W^2 \cdot W^1 = W_{eff} \longrightarrow y = W_{eff} x$

  (single neuron!)

- For which activation functions is a deep approximately linear when $x \sim N(0, \sigma^2)$ with $\sigma^2 \ll 1$?

1) Sigmoid   $s(x) = \dfrac{1}{1 + 2^{-x}}$

   Taylor expansion in $\varphi$:   $s(x) \approx \dfrac{1}{2} + \dfrac{x}{4} - \dfrac{x^3}{48} + \ldots$

   For small $\sigma$   $s(x) \approx \dfrac{1}{2} + \dfrac{x}{4}$   Not linear: constant bias breaks linearity over multiple layers ✗

2) $\tanh(x)$

   Taylor expansion   $s(x) = x - \dfrac{x^3}{3} + \dfrac{2x^5}{15} - \ldots$

   For small $\sigma$   $s(x) = x$   Linear ✓

3) $\text{Relu}(x) = \max(0, x)$

   $x < 0 \longrightarrow$ output 0
   $x > 0 \longrightarrow$ output $x$

   ↳ Normal distribution around $\emptyset$ means that whole interval are mapped into $\emptyset$

   Non linearity !

a) $\text{Selu}(x) = \begin{cases} \lambda x & x > 0 \\ \lambda \alpha (e^x - 1) & x \leq 0 \end{cases}$

Taylor expansion $\begin{cases} \lambda x & x > 0 \\ \lambda \alpha (x + \frac{x^2}{2} + \cdots) \approx \lambda \alpha x & x \leq 0 \end{cases}$

$\hookrightarrow$ this is piecewise linear (but not globally) slope differs between $x > 0$ and $x < 0$