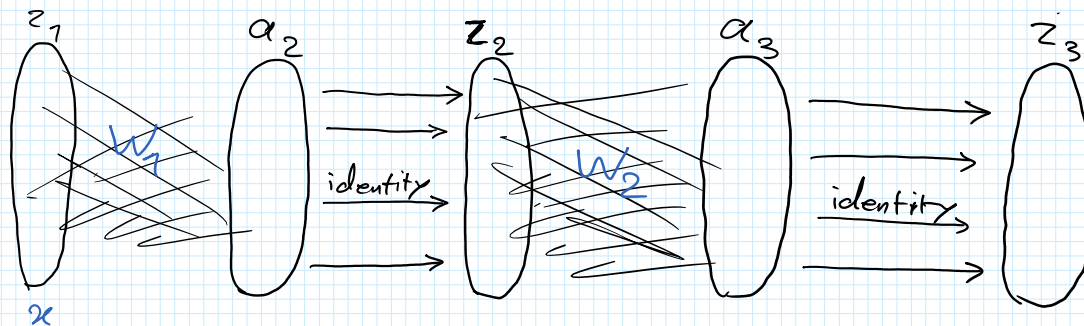


1 Activation functions

Problem 1: Why do we use non-linear activation functions in neural networks? What purpose do they serve?



$$a_2 = W_1^T x$$

$$z_2 = a_2 = W_1^T x$$

$$a_3 = W_2^T z_2 = W_2^T a_2 = W_2^T W_1^T x$$

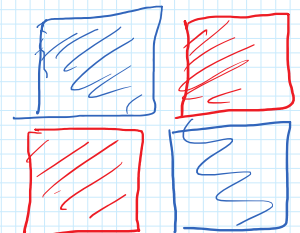
linear transformation
output
input

$$z_3 = a_3 = W_2^T W_1^T x$$

What if we had more than 2 layers?

linear transformation

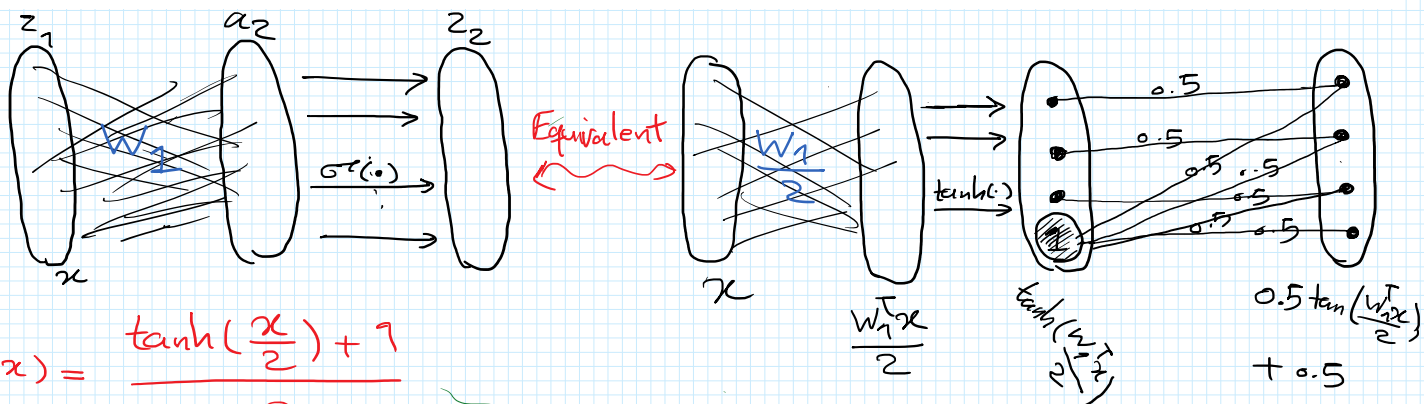
$$z_L = W_{L-1}^T \dots W_2^T W_1^T x$$



What is the problem with it?

It fails, e.g., on XOR dataset.

Problem 2: Consider a neural network where the hidden units use the sigmoid activation function. Show that there exists an equivalent network, which computes exactly the same function, but with hidden units using $\tanh(x)$ as activation functions.

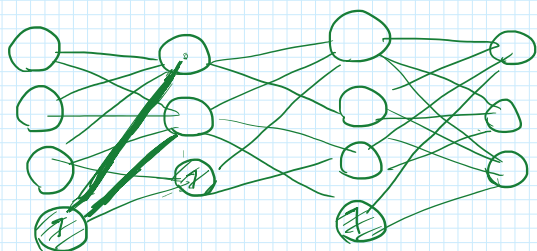


Let's derive:

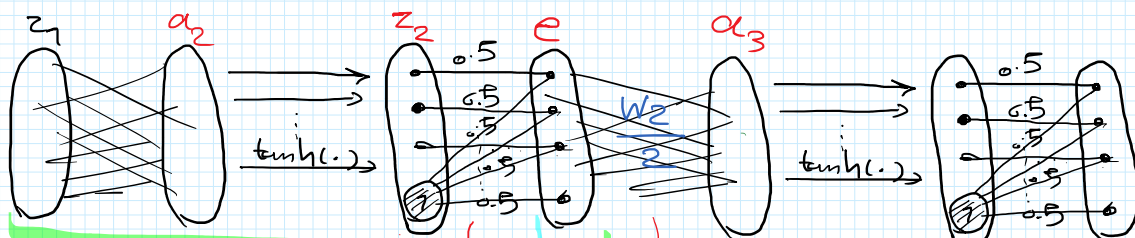
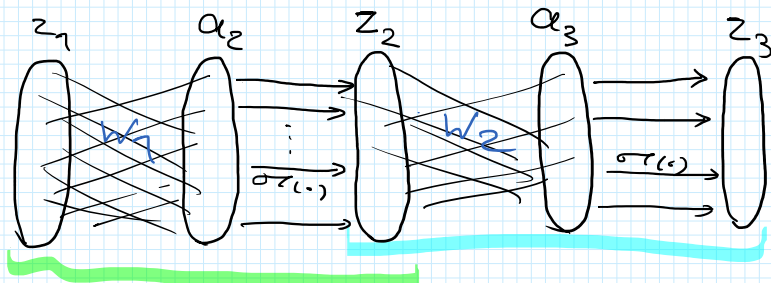
$$\frac{1}{2} \left(\frac{e^{0.5x} - e^{-0.5x}}{e^{0.5x} + e^{-0.5x}} + 1 \right)$$

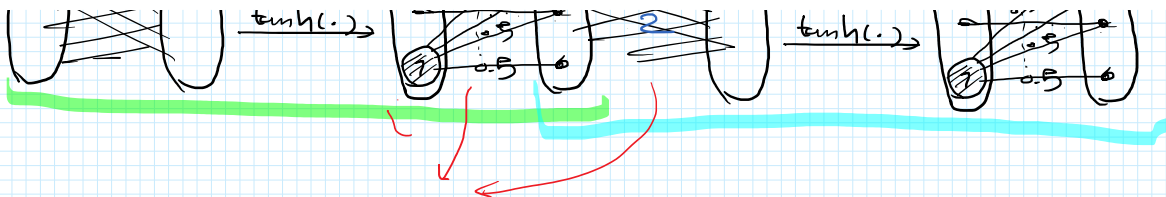
$$= \frac{1}{2} \left(\frac{e^{0.5x} - e^{-0.5x} + e^{0.5x} + e^{-0.5x}}{e^{0.5x} + e^{-0.5x}} \right)$$

$$= \frac{1}{2} \left(\frac{2e^{0.5x}}{e^{0.5x} + e^{-0.5x}} \right) = \left(\frac{1}{1 + e^{-x}} \right)$$



What if we had more than 1 layer?





combine these two.

$$e = \underbrace{\begin{bmatrix} 0.5 & 0.6 & 0.5 \\ 0.5 & 0.5 & 0.5 \end{bmatrix}}_A z_2$$

$$a_3 = \frac{w_2^T}{2} e$$

$$a_3 = \left(\frac{w_2^T}{2} \quad A^T \right) z_2 \rightarrow \text{A new linear layer}$$

Problem 3: We already know that the derivative of the sigmoid activation function can be expressed in terms of the function value itself. Show that the derivative of the tanh activation function can also be expressed in terms of the function value itself. Why is this a useful property?

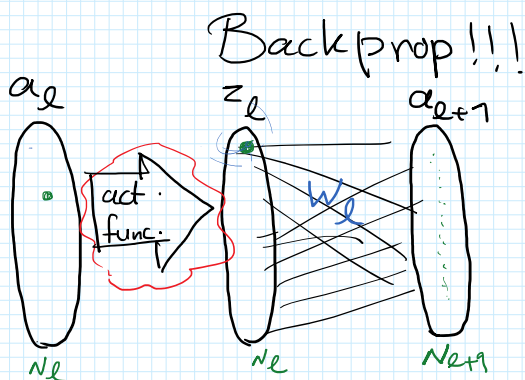
$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Task: $\tanh'(x) = 1 - \tanh^2(x)$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\begin{aligned} \tanh'(x) &= \frac{(e^x + e^{-x})(e^x - e^{-x}) - (e^x - e^{-x})(e^x + e^{-x})}{(e^x + e^{-x})^2} \\ &= \frac{(e^x + e^{-x})^2}{(e^x + e^{-x})^2} - \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right)^2 = \boxed{1 - \tanh^2(x)} \end{aligned}$$

Why do we need derivative of activation function?



We have: $\frac{\partial E_n}{\partial a_{l+1}} = \delta_{l+1}^{(n)}$

We want: $\frac{\partial E_n}{\partial a_l} = \delta_l^{(n)}$

$$\frac{\partial E_n}{\partial a_l} = \frac{\partial E_n}{\partial a_{l+1}} \frac{\partial a_{l+1}}{\partial z_l} \frac{\partial z_l}{\partial a_l}$$

Dimensions: $1 \times N_{l+1}$, $[N_{l+1} \times N_l]$, $[N_l \times N_l]$

$$\left\{ \frac{\partial z_{l+1}}{\partial a_{l+2}} = 0 \right.$$

if activation function is $\tanh(\cdot)$

$$\Rightarrow \frac{\partial z_l}{\partial a_l} = \begin{bmatrix} \tanh'(a_{l1}) & \tanh'(a_{l2}) \\ \tanh'(a_{l1}) & \tanh'(a_{l2}) \end{bmatrix}$$

Problem 4: In machine learning you quite often come across problems which contain the following quantity

$$y = \log \sum_{i=1}^N e^{x_i}$$

For example if we want to calculate the log-likelihood of a neural network with a softmax output we get this quantity due to the normalization constant. If you try to calculate it naively, you will quickly encounter underflows or overflows, depending on the scale of x_i . Despite working in log-space, the limited precision of computers is not enough and the result will be ∞ or $-\infty$.

To combat this issue we typically use the following identity:

$$y = \log \sum_{i=1}^N e^{x_i} = a + \log \sum_{i=1}^N e^{x_i - a}$$

for an arbitrary a . This means, you can shift the center of the exponential sum. A typical value is setting a to the maximum ($a = \max_i x_i$), which forces the greatest value to be zero and even if the other values would underflow, you get a reasonable result.

Your task is to show that the identity holds.

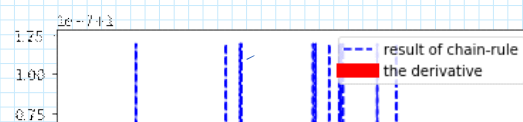
example: $\log(e^{-100} + e^{-101} + e^{-102})$

Naive: $t1 = \text{np.exp}(-100)$
 $t2 = \text{np.exp}(-101)$
 $t3 = \text{np.exp}(-102)$
 $\text{result} = \text{np.log}(t1 + t2 + t3)$
nan

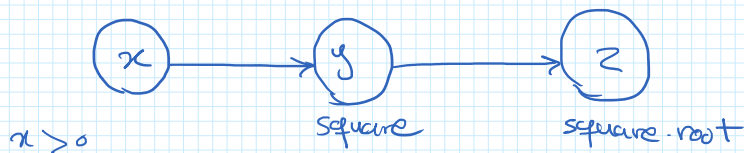
$\log(e^{-100} (1 + e^{-1} + e^{-2}))$
 $= -100 + \log(1 + e^{-1} + e^{-2})$
✓ 0 ≠

$$\begin{aligned} \log(e^{x_1} + e^{x_2} + \dots + e^{x_n}) &= \log(e^a (e^{x_1-a} + \dots + e^{x_n-a})) \\ &= \log(e^a) + \log(e^{x_1-a} + \dots + e^{x_n-a}) \\ &= a + \log(e^{x_1-a} + \dots + e^{x_n-a}) \end{aligned}$$

Example: Numerical instability



Example: Numerical instability



$$y = x^2 \quad z = y^{\frac{1}{2}}$$

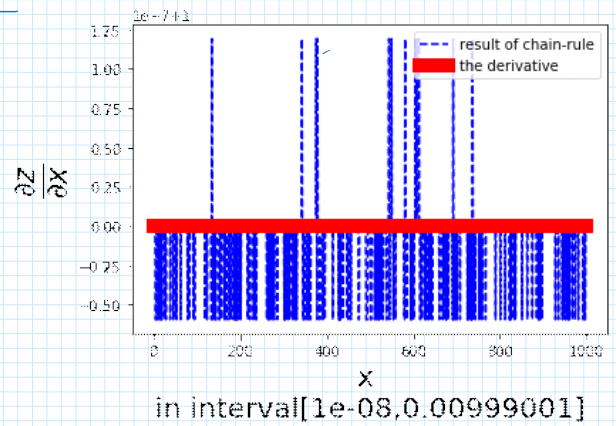
$$\frac{\partial z}{\partial x} = \frac{\partial}{\partial x} x = 1$$

Using chain-rule: $\left. \frac{\partial z}{\partial x} \right|_{x_0}$

Forward Pass: $y_0 = x_0^2 \quad z_0 = \sqrt{y_0}$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} = \left[\frac{1}{2} y^{-\frac{1}{2}} \right] [2x]$$

$$\Rightarrow \left. \frac{\partial z}{\partial x} \right|_{x_0} = y_0^{-\frac{1}{2}} x_0$$



if $x_0 \rightarrow 0$

Problem 5: Similar to the previous exercise we can compute the output of the softmax function $\pi_i = \frac{e^{x_i}}{\sum_{i=1}^N e^{x_i}}$ in a numerically stable way by introducing an arbitrary a :

$$\frac{e^{x_i}}{\sum_{i=1}^N e^{x_i}} = \frac{e^{x_i-a}}{\sum_{i=1}^N e^{x_i-a}}$$

often chosen $a = \max_i x_i$. Show that the above identity holds.

$$\text{Softmax}(x_1, \dots, x_N) = \left[\frac{e^{x_1}}{\sum_i e^{x_i}}, \frac{e^{x_2}}{\sum_i e^{x_i}}, \dots, \frac{e^{x_N}}{\sum_i e^{x_i}} \right]$$

Example: $\frac{e^{-100}}{e^{-100} + e^{-101} + e^{-102}} = \frac{0}{0} = \text{NaN}$

\nearrow small \nearrow small

The Better Way: $\frac{e^{-100} \cdot e^0}{e^{-100} \cdot (e^0 + e^{-1} + e^{-2})} = \frac{1}{e^0 + e^{-1} + e^{-2}} \checkmark$

$$\frac{e^{x_i-a}}{\sum_j e^{x_j-a}} = \frac{\cancel{e^{-a}} e^{x_i}}{\cancel{e^{-a}} \left(\sum_j e^{x_j} \right)} = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

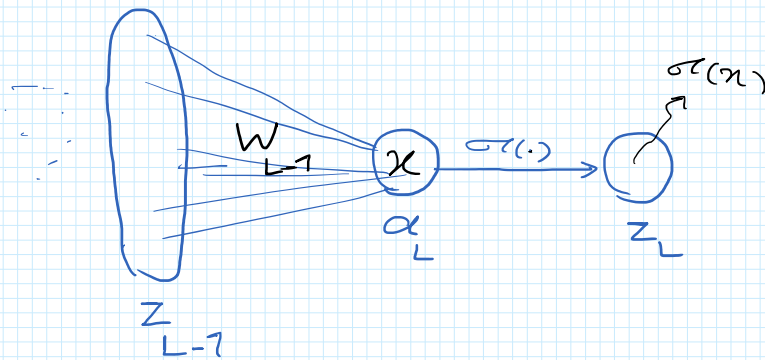
Problem 6: Let the logits (the values before applying sigmoid) of a single training example be x and the corresponding label be y . The sigmoid logistic loss (i.e. binary cross-entropy) for this example is then:

$$-(y \log(\sigma(x)) + (1 - y) \log(1 - \sigma(x)))$$

To ensure stability and avoid overflow, usually implementations use this equivalent formulation:

$$\max(x, 0) - x \cdot y + \log(1 + e^{-\text{abs}(x)})$$

Your task is to show that the equivalence holds.



$$\Rightarrow \text{binary cross-entropy: } -[y \log(\sigma(x)) + (1-y) \log(1-\sigma(x))]$$

$$\text{cost} = -\left[y \log\left(\frac{1}{1+e^{-x}}\right) + (1-y) \log\left(1 - \frac{1}{1+e^{-x}}\right) \right]$$

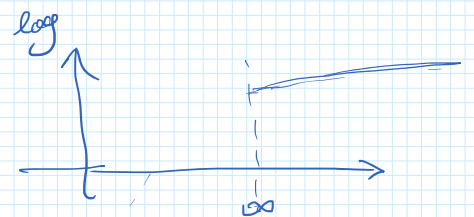
$$= -\left[-y \log(1+e^{-x}) + (1-y) \log\left(\frac{e^{-x}}{1+e^{-x}}\right) \right]$$

$$= -\left[-y \log(1+e^{-x}) + (1-y)(-x - \log(1+e^{-x})) \right]$$

$$\text{cost} = x - xy + \log(1+e^{-x})$$

if $x > 0 \Rightarrow e^{-x}$ small $\Rightarrow \checkmark$

$-x$...



if $x > 0 \Rightarrow e^{-x}$ small $\Rightarrow \checkmark$

if $x < 0 \Rightarrow e^{-x}$ might be too large

→ We can rewrite it as:

$$\text{cost} = x - xy + \log(e^{-x}(e^x + 1))$$

$$= \cancel{x} - xy + \log(\cancel{e^{-x}}) + \log(e^x + 1)$$

$$= -xy + \log(e^x + 1) \Rightarrow x < 0$$

$$\Rightarrow \text{cost} = \max(0, x) - xy + \log(1 + e^{-|x|})$$