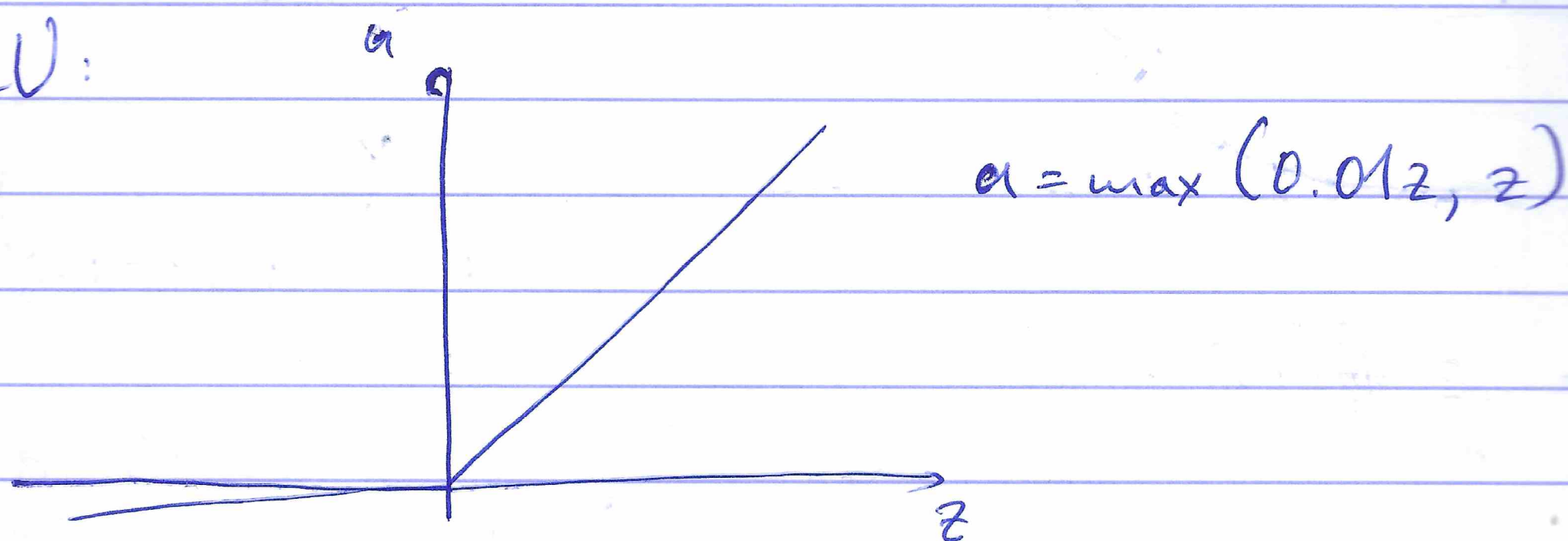


For a lot of the space of z , the slope of a is very different from 0, so the net will learn much faster.

Leaky ReLU:



Why do you need non-linear activation functions?

Say $g(z) = z$ (linear activ. func.):

$$\left. \begin{aligned} z^{[1]} &= W^{[1]}x + b^{[1]} \\ a^{[1]} &= z^{[1]} \\ z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} &= z^{[2]} \end{aligned} \right\}$$

$$\begin{aligned} a^{[1]} &= z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[2]} &= z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} &= W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]} \\ &= (W^{[2]}W^{[1]})x + (W^{[2]}b^{[1]} + b^{[2]}) \\ &= W'x + b' \end{aligned}$$

A linear hidden layer is practically useless. The combo of 2 linear functions is a non-linear function.

This might only be useful if you're doing ML on a regression problem. ^{maybe} at the output layer. If you're predicting \$0 ... \$1,000,000 (housing prices?)