

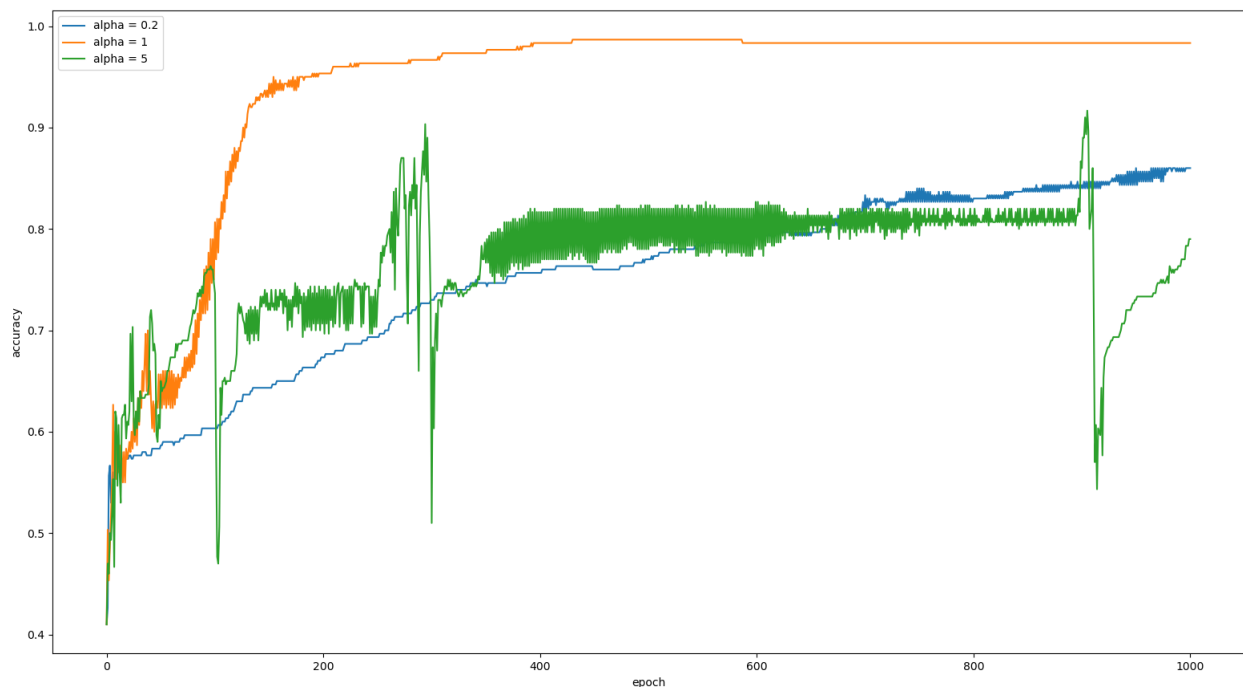
Report

Computer lab. Backpropagation

Morhunenko Mykola.

Assignment 1.

The convergence plots for three settings of the learning rate $\alpha \in \{0.2, 1, 5\}$ are:



The result convergence improves with increasing the learning rate goes up to some point, and after that, it decreases and shows much worse results. It is because the learning rate is proportional to the chosen method parameter - Gradient Descent - and up to some point, decreasing of the parameter alpha improves the convergence, but after some point it will oscillate and show much worse results. The optimization function results are moving around the optimal value but can not reach it because of “overshooting”, which decreases the accuracy.

Assignment 2.

1) softmax:

$$p_k(s) = \frac{e^{s_k}}{\sum_{c=1}^K e^{s_c}}$$

2) cross-entropy loss:

$$L = - \sum_{k=1}^K t_k \log(p_k)$$

Forward and backward messages for a compound layers 1) and 2):

forward:

$$L_p = - \sum_{k=1}^K t_k (s_k - \log(\sum_{c=1}^K e^{s_c}))$$

backward:

$$\frac{\partial L}{\partial s_k} = t_k - p_k$$

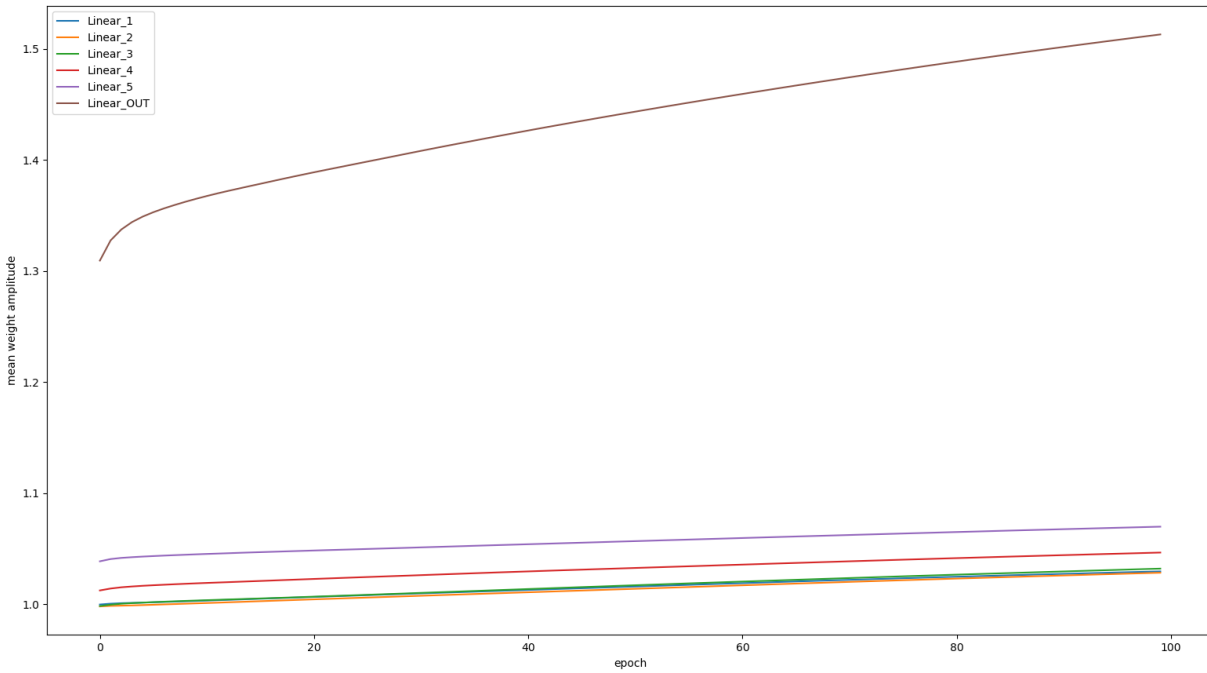
$s = (s_1, \dots, s_K)$
 $p = (p_1, \dots, p_K)$
 $t = (t_1, \dots, t_K)$ - one-hot encoded vector of targets.

\equiv Show that softmax is invariant to shift in inputs, $p_k(s'_k) = p_k(s_k)$, $s'_k = s_k + d$ for $k \in \{1, \dots, K\}$, $d \in \mathbb{R}$.

$$p_k(s'_k) = \frac{e^{s_k + d}}{\sum_{c=1}^K e^{s_c + d}} = \frac{e^{s_k} e^d}{\sum_{c=1}^K e^{s_c} e^d} = \frac{e^d e^{s_k}}{e^d \sum_{c=1}^K e^{s_c}} = p_k(s_k) \blacktriangle$$

Assignment 3.

The plot shows the development of mean weight amplitude for each linear layer over the training epoch normalized w.r.t. The initial mean amplitude is shown below:



The weights are initialized randomly, so the training loss is higher for the first epochs. Because of the specifics of the Gradient Descent method, the weights amplitude firstly increases fast at the beginning and then slows down, because the training loss decreases.