

Team AIAufgabe 1:

1a) Sei $X_{ki} \sim \mathcal{N}(0, 1)$ jeweils die ZV, welche w_{ki} beschreibt.

$$\Rightarrow \sum_{k=1}^m x_k X_{ki} + \underbrace{b_k}_{=0} = \sum_{k=1}^{\lfloor \frac{m}{2} \rfloor} X_{ki} \sim \mathcal{N}\left(\sum_{k=1}^{\lfloor \frac{m}{2} \rfloor} 0, \sqrt{\sum_{k=1}^{\lfloor \frac{m}{2} \rfloor} 1}\right) \\ = \mathcal{N}\left(0, \sqrt{\lfloor \frac{m}{2} \rfloor}\right) = \mathcal{N}(0, \sqrt{150})$$

b) Durch die breite Varianz entstehen eventuell dendritische Potential weit weg von 0. ~~Siehe~~ Dadurch würde das Problem von vanishing gradient auftauchen.

$$c) \sum_{k=1}^m x_k X_{ki} + b_k = \sum_{k=1}^{\lfloor \frac{m}{2} \rfloor} X_{ki} \sim \mathcal{N}\left(0, \left(\sum_{k=1}^{\lfloor \frac{m}{2} \rfloor} \frac{1}{n_{in}}\right)^{\frac{1}{2}}\right) = \mathcal{N}\left(0, \left(\frac{m}{2n_{in}}\right)^{\frac{1}{2}}\right) \\ = \mathcal{N}\left(0, \left(\frac{1}{2}\right)^{\frac{1}{2}}\right)$$

d) Die W'keit, dass die dendritischen Potentiale jetzt hohe Werte annehmen, ist deutlich geringer.

2.) I entspricht A und II entspricht B.

Denn der tanh staucht die breite Verteilung, die durch I verursacht wird, an den Rändern. Wodurch bei A viele ~ 1 und ~ -1 entstehen. Außerdem ist die Ableitungen an diesen Stellen fast 0, was zur Abbildung der Gradienten passt. Die Abbildung passt auch zu den berechneten Verteilungen

3) Die ~~Xavier-Initialisierung~~ skalierte Normalverteilung reicht das Problem der Standardnormalverteilung einfach nur eine Schicht weiter.

Die Xavier-Initialisierung löst dieses Problem und pro Schicht wird dadurch unter Beachtung aller Neuronen, die Varianz beibehalten.

Aufgabe 2:

1a) Zunächst:
$$\frac{\partial}{\partial w} \frac{\lambda}{2} \sum_{\ell=1}^L \|w^\ell\|_F^2 = \frac{\lambda}{2} 2w = \lambda \cdot w$$

$$\begin{aligned} \Rightarrow w(t+1) &= w(t) - \eta \nabla E(w(t)) \\ &= w(t) - \eta \nabla E_0(w(t)) - \eta \frac{\partial}{\partial w} \frac{\lambda}{2} \sum_{\ell=1}^L \|w^\ell\|_F^2 \\ &= w(t) - \eta \nabla E_0(w(t)) - \eta \lambda w(t) \\ &= (1 - \lambda \eta) w(t) - \eta \nabla E_0(w(t)) \end{aligned}$$

b) i) Mit den Werten eingesetzt ergibt sich:

$$\begin{aligned} w(t+1) &= w(t) (1 - 0,8 \cdot 0,5) - 0 \\ &= 0,6 \cdot w(t) \end{aligned}$$

$$\Rightarrow w(t) = (0,6)^t \cdot 2 \quad \Rightarrow w(10) \approx 0,042$$

ii) Das Gewicht sinkt logarithmisch.

iii) $\lim_{t \rightarrow \infty} w(t) = \lim_{t \rightarrow \infty} (0,6)^t \cdot 2 = 0$ (da streng monoton fallend und nach unten beschränkt)

iv) In diesem Fall wird jede Iteration durch die Ableitung beeinflusst. Also wird das Gewicht immer verändert, solange es nicht optimal ist.

c) Regularisierung hilft vanishing gradients einzudämmen, da jedes Gewicht natürlich ~~verändert~~ verkleinert wird (absolut) und nach einigen Iterationen sollten die meisten Gewichte nahe um 0 gestreut sein.

2.) Wie man auf Blatt 3 sieht, fällt der Faktor des Inputs in der Ableitung nach b weg. D.h. der Input hat deutlich geringeren Einfluss auf die Anpassung des Bias und daher auch verzerrte Inputs.

3.)

