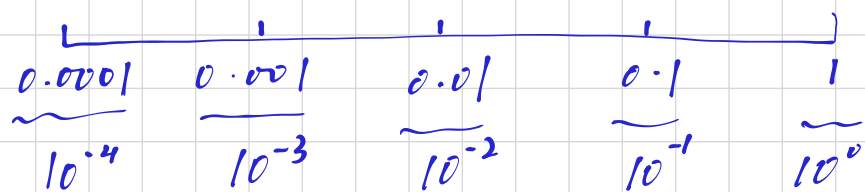


Empire DNN - Week 3

1. Hyperparameter Tuning

(1) Appropriate scale for hyperparameters



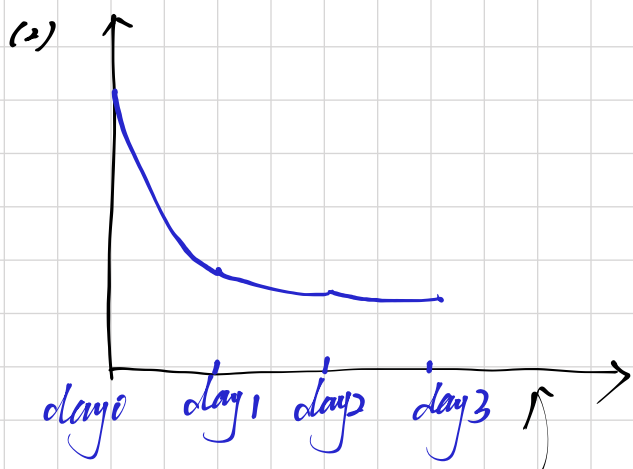
Hyperparameters for exponentially weighted averages

$$\Rightarrow \gamma \in [-3, -1]$$

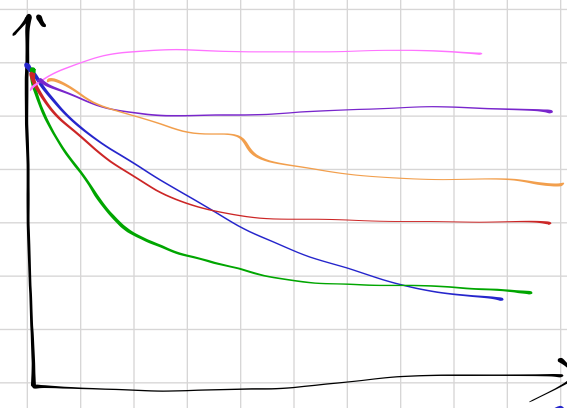
$$1 - \beta = 10^\gamma$$

Sample more densely when $\beta \rightarrow 1$

$$\beta = 1 - 10^\gamma \rightarrow \beta = 0.9, 0.99, 0.999$$



Boostrapping one model \leftarrow parallel



Training many models in parallel

2. Batch Normalization

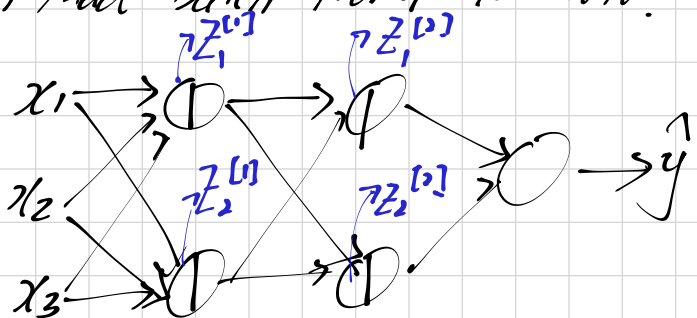
(1) Give some intermediate value in NN: $z^{(1)}, \dots, z^{(m)}$

$$\begin{cases} \mu = \frac{1}{m} \sum_i z^{(i)} \\ \sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2 \\ z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \\ \tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta \end{cases}$$

$$\begin{aligned} \gamma &= \sqrt{\sigma^2 + \epsilon} \\ \beta &= \mu \\ \text{then } \tilde{z}^{(i)} &= z^{(i)} \end{aligned}$$

\Rightarrow use $\tilde{z}^{(i)}$ instead of $z^{(i)}$

(2) Add Batch Norm to NN:



$$\begin{aligned} x &\xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{(BN)}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \rightarrow a^{[1]} = f(\tilde{z}^{[1]}) \\ &\quad \vdots \leftarrow a^{[2]} \leftarrow \tilde{z}^{[2]} \xrightarrow[\text{(BN)}]{\beta^{[2]}, \gamma^{[2]}} z^{[2]} \xrightarrow{w^{[2]}, b^{[2]}} \end{aligned}$$

For mini-batches:

$$\begin{cases} X^{[1]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \rightarrow g^{[1]}(\tilde{z}^{[1]}) = a^{[1]} \rightarrow \dots \\ X^{[2]} \rightarrow \dots \\ X^{[3]} \rightarrow \dots \end{cases}$$

parameters: $w^{[1]}, b^{[1]}, \beta^{[1]}, \gamma^{[1]}$

Batch Norm as regularization:

① Each mini-batch is scaled by the mean/variance computed on just that mini-batch.

② This adds some noise to the values $z^{[1]}$ within that mini-batch. So similar to dropout, it adds some noise to each hidden layer's activations.

③ This has a slight regularization effect.

3. Multi-class classification

(1) Softmax Regression

$$x \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Layer L

$$z^{[1]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \rightarrow \sum_{j=1}^4 t_j = 176.3$$

$$a^{[1]} = \frac{t_i}{\sum_{j=1}^4 t_j}, \quad a_i^{[1]} = \frac{t_i}{\sum_{j=1}^4 t_j}$$

$$\begin{aligned} \frac{e^5}{176.3} &= 0.842 \\ \frac{e^2}{176.3} &= 0.042 \\ \frac{e^{-1}}{176.3} &= 0.02 \\ \frac{e^3}{176.3} &= 0.114 \end{aligned}$$

Softmax regression generalizes logistic regression to C classes.