

Deep Learning - HW1

Afek Adler & Or Wolkomir

November 25, 2019

1 Theory

1.1 Question 1 - network architecture

- (a) The input shape is $(m, 10)$
- (b) W_h shape is $(10, 50)$ and b_h shape is $(, 50)$
- (c) W_o shape is $(50, 3)$ and b_o shape is $(, 3)$
- (d) y shape is $(m, 3)$
- (e) W_h shape is $(50, 3)$ and b_h shape is $(, 3)$

1.2 Question 2 - # of parameters

- (a) **Conv layer 1**
Kernel = 3
Channels = 3
Output = 100
Bias := Output
Therefore:

$$\begin{aligned}\#_of_parameters &= Kernel^2 * Channels * Output + Bias \Rightarrow \\ &3^2 * 100 + 100 = 2,800\end{aligned}$$

- (b) **Conv layer 2**
Kernel = 3
Channels = 100
Output = 200
Bias := Output
Therefore:

$$\begin{aligned}\#_of_parameters &= Kernel^2 * Channels * Output + Bias \Rightarrow \\ &3^2 * 100 * 200 + 200 = 180,200\end{aligned}$$

(c) **Conv layer3**

Kernel = 3

Channels = 200

Output = 400

Bias := Output

Therefore:

$$\begin{aligned}\#_of_parameters &= Kernel^2 * Channels * Output + Bias \Rightarrow \\ 3^2 * 200 * 400 + 400 &= 720,400\end{aligned}$$

(d) **Total**

$$total_ \#_of_parameters = 720,400 + 180,200 + 2,800 = 903,400$$

1.3 Question 3 - deriving BatchNorm gradients

(a) $\frac{\partial f}{\partial \gamma} =$

$$\frac{\partial f}{\partial y} \cdot \frac{\partial y}{\partial \gamma} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \cdot \hat{x}_i \quad (1)$$

(b) $\frac{\partial f}{\partial \beta} =$

$$\frac{\partial f}{\partial y} \cdot \frac{\partial y}{\partial \beta} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \frac{\partial y}{\partial \beta} = \sum_{i=1}^m \frac{\partial f}{\partial y_i} \quad (2)$$

(c) $\frac{\partial f}{\partial \hat{x}_i} =$

$$\frac{\partial f}{\partial y} \cdot \frac{\partial y}{\partial \hat{x}_i} = \frac{\partial f}{\partial y_i} \cdot \gamma \quad (3)$$

(d) $\frac{\partial f}{\partial \sigma^2} =$

$$\frac{\partial f}{\partial \hat{x}} \cdot \frac{\partial \hat{x}}{\partial \sigma^2} = \frac{\partial y}{\partial \hat{x}_i} = \frac{\partial f}{\partial y_i} \cdot \gamma \quad (4)$$

(e) $\frac{\partial f}{\partial \mu} =$

2 Practical

We wanted to compare apples to apples so we tried to keep most hyper-parameters consistent across experiments, for educational purposes (it is clear that if we would do grid search for example most likely that hyper-parameters will vary between experiments). As well, we would like to note that we did not perform hyper-parameter tuning, because the networks achieved sufficient performance out of the box. This is way we also did not use a validation data set.

As we wanted to compare apples to apples we regularized the model only in one layer, after the first fully connected layer. We could have also done differently (another hyper-parameter). We used the following hyper-parameters across all the experiments :

- (a) Optimizer - ADAM with default β_1 and β_2
- (b) Learning rate - 0.001
- (c) Epochs - 10. We would like to note that we set seed for each experiment, so if it was a real model (in production), we can perform early stopping and choose the best score. But than it's better to do have a test score as well (because we are "playing" with the model) in order to get a real estimate about the algorithm accuracy.

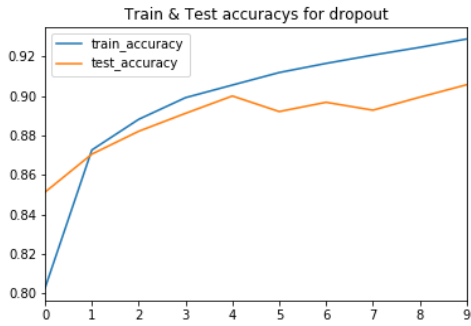
Hyper-parameters that are unique for each experiment -

1. BatchNorm with hyper-parameter affine = False
2. Dropout (with 10
3. Weight decay (for all weights, not only one layer) - $2 \cdot 10^{-4}$
4. No regularization

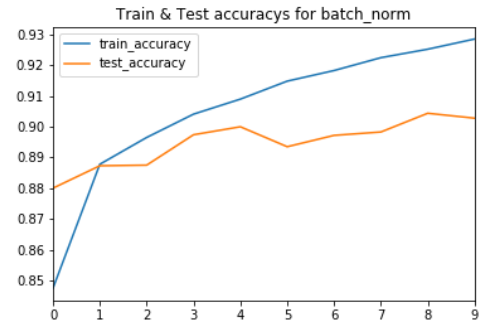
Conclusions

- (a) For all the experiments after 2-3 epochs the model over-fits.
- (b) Even though over-fitting happens validation scores keep increasing with the epochs.
- (c) All the models are pretty much the same (around 90% accuracy). So in order to argue which model is better it's better to perform an hypothesis test.

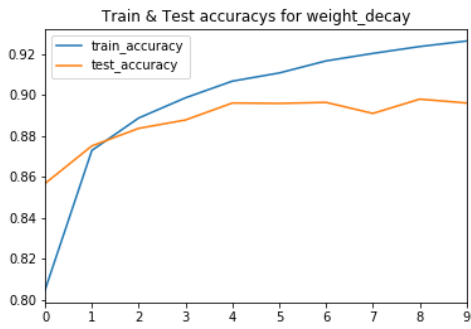
2.1 Graphs



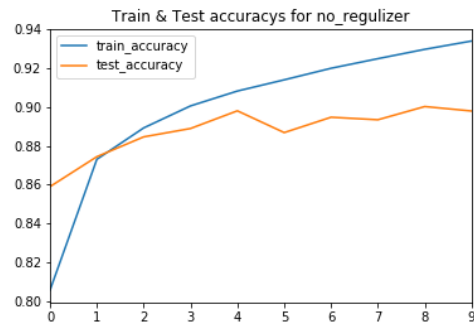
(a)



(b)



(c) 3



(d)

Figure 1: Convergence graphs