# CSE 291 I00: Machine Learning for 3D Data
# Homework 1

Qiaojun Feng
A53235331
Department of Electrical and Computer Engineering
University of California, San Diego
qjfeng@ucsd.edu

## 1

Compute the derivative of the cross-entropy loss $L$ w.r.t $o$, i.e compute $\frac{\partial L}{\partial o}$.

$$\frac{\partial L}{\partial o} = \frac{\partial L}{\partial p}\frac{\partial p}{\partial o}$$
$$L = -\sum_{\text{for all class } c} y_c log(p_c)$$
$$\frac{\partial L}{\partial p_c} = -\frac{y_c}{p_c}$$
$$p_c = \frac{e^{o_c}}{\sum_{j=0}^{K} e^{o_j}}$$
$$\frac{\partial p_c}{\partial o_i} = \begin{cases} p_c(1-p_i) & i = c \\ -p_c p_i & i \neq c \end{cases}$$
$$\frac{\partial L}{\partial o_i} = \sum_{\text{for all class } c} \frac{\partial L}{\partial p_c}\frac{\partial p_c}{\partial o_i}$$
$$= -y_i(1-p_i) + \sum_{c \neq i} y_c p_i$$
$$= -y_i + (y_i + \sum_{c \neq i} y_c)p_i$$
$$= p_i - y_i$$

The detailed derivation of softmax's derivative:
For $i = c$:

$$\frac{\partial p_c}{\partial o_i} = \frac{e^{o_c}\sum_{j=0}^{K} e^{o_j} - e^{o_c}e^{o_i}}{(\sum_{j=0}^{K} e^{o_j})^2}$$
$$= \frac{e^{o_c}}{\sum_{j=0}^{K} e^{o_j}} \times \frac{\sum_{j=0}^{K} e^{o_j} - e^{o_i}}{\sum_{j=0}^{K} e^{o_j}}$$
$$= p_c(1-p_i)$$

For $i \neq c$:

$$\frac{\partial p_c}{\partial o_i} = \frac{-e^{o_c}e^{o_i}}{(\sum_{j=0}^{K} e^{o_j})^2}$$
$$= -\frac{e^{o_c}}{\sum_{j=0}^{K} e^{o_j}} \times \frac{e^{o_i}}{\sum_{j=0}^{K} e^{o_j}}$$
$$= -p_c p_i$$

## 2

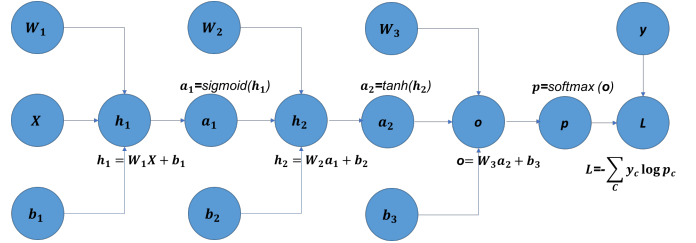(1)Draw the computational graph of the network.



Fig. 1.  computational graph of network

(2)Perform Backpropagation to compute the gradients of loss $L$ w.r.t all the weights and biases,
i.e compute $\frac{\partial L}{\partial W_1}, \frac{\partial L}{\partial W_2}, \frac{\partial L}{\partial W_3}, \frac{\partial L}{\partial b_1}, \frac{\partial L}{\partial b_2}, \frac{\partial L}{\partial b_3}$.

$$\frac{\partial L}{\partial W_3} = \sum_{i=0}^{K} \frac{\partial L}{\partial o_i}\frac{\partial o_i}{\partial W_3}$$
$$= \sum_{i=0}^{K} (p_i - y_i)a_2^T$$
$$\frac{\partial L}{\partial b_3} = \sum_{i=0}^{K} \frac{\partial L}{\partial o_i}\frac{\partial o_i}{\partial b_3}$$
$$= \sum_{i=0}^{K} (p_i - y_i)$$
$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial a_2}\frac{\partial a_2}{\partial h_2}\frac{\partial h_2}{\partial W_2}$$
$$= W_3^T \sum_{i=0}^{K} (p_i - y_i)(1 - a_2^2)a_1^T$$
$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial a_2}\frac{\partial a_2}{\partial h_2}\frac{\partial h_2}{\partial b_2}$$
$$= W_3^T \sum_{i=0}^{K} (p_i - y_i)(1 - a_2^2)$$
$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial a_1}\frac{\partial a_1}{\partial h_1}\frac{\partial h_1}{\partial W_1}$$
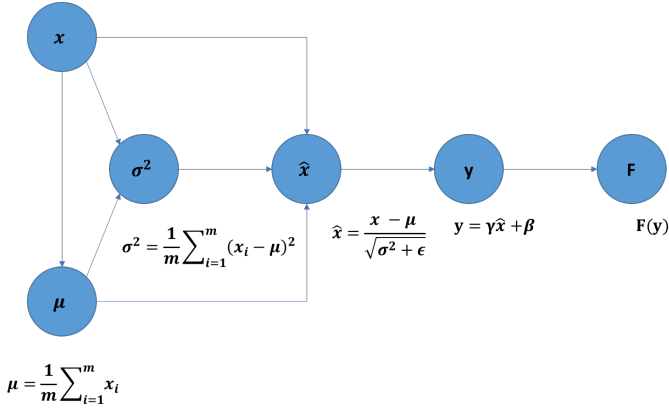$$= W_2^T W_3^T \sum_{i=0}^{K} (p_i - y_i)(1 - a_2^2)a_1(1 - a_1)X^T$$

Fig. 2. computational graph of BatchNorm layer

$$\mu = \frac{1}{m}\sum_{i=1}^{m} x_i$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial a_1}\frac{\partial a_1}{\partial h_1}\frac{\partial h_1}{\partial b_1}$$

$$= W_2^T W_3^T \sum_{i=0}^{K}(p_i - y_i)(1 - a_2^2)a_1(1 - a_1)$$

3

(1)Draw the computational graph of the batchnorm layer. (2)Consider a function $F(y_i)$ and compute the derivatives of $F(y_i)$ w.r.t $\gamma$, $\beta$ and $x$, i.e compute $\frac{\partial F}{\partial \gamma}, \frac{\partial F}{\partial \beta}, \frac{\partial F}{\partial x}$. Assume $\frac{\partial F}{\partial y_i}$ is given.

If we want derive the derivative from the batch of data. That is, function $F(y_1, \ldots, y_m)$. Then we should have:

$$\frac{\partial F}{\partial \gamma} = \sum_{i=1}^{m}\frac{\partial F}{\partial y_i}\frac{\partial y_i}{\partial \gamma}$$

$$= \sum_{i=1}^{m}\frac{\partial F}{\partial y_i}\hat{x}_i$$

$$\frac{\partial F}{\partial \beta} = \sum_{i=1}^{m}\frac{\partial F}{\partial y_i}\frac{\partial y_i}{\partial \beta}$$

$$= \sum_{i=1}^{m}\frac{\partial F}{\partial y_i}$$

$$\frac{\partial F}{\partial \hat{x}_i} = \frac{\partial F}{\partial y_i}\gamma$$

$$\frac{\partial F}{\partial \sigma_{\mathcal{B}}^2} = \sum_{i=1}^{m}\frac{\partial F}{\partial \hat{x}_i}(x_i - \mu_{\mathcal{B}})\frac{-1}{2}(\sigma_{\mathcal{B}}^2 + \epsilon)^{-2/3}$$

$$\frac{\partial F}{\partial \mu_{\mathcal{B}}} = \sum_{i=1}^{m}\frac{\partial F}{\partial \hat{x}_i}\frac{-1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$\frac{\partial F}{\partial x_i} = \frac{\partial F}{\partial \hat{x}_i}\frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \frac{\partial F}{\partial \sigma_{\mathcal{B}}^2}\frac{2(x_i - \mu_{\mathcal{B}})}{m} + \frac{\partial F}{\partial \mu_{\mathcal{B}}}\frac{1}{m}$$

4

I referred to Udacity's github deep-learning repo and built the network on its fundamental.

(1)(2)Construct a model using fully connected layers and ReLu layers to solve mnist classification problem. Plot a graph demonstrating how the loss function decreases over the number of iterations. Compare with and without batch normalization layers. Compare different optimization methods.
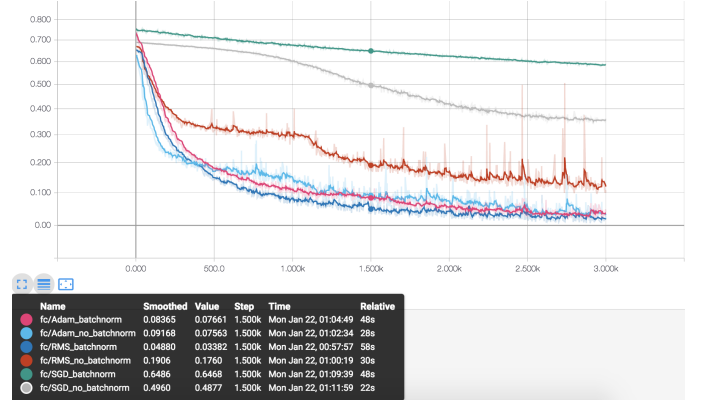


Fig. 3. loss of fully connected network

For the dataset, there are 60000 as train set and 10000 as test set. We split the train set and use 5000 of them as validation set. So the remaining train set has 55000 images. Use batch size 110, we need 500 batches to go through the whole train set.

We designed a 8-fully-connected-layer network. The numbers of neurons in each layer are 256,128,64,32,16,8,4,10. Also we tried stochastic gradient descent, RMSProp and Adam optimizer. We can conclude that RMSProp and Adam are better than SGD, while batch normalization helps all of them.

Batch normalization helps because it normalize the input to each layer and make the model training easier. The weight will be changed in a smaller range without broadly search. Also since on each layer the derivative will be normalized, the back-propagation is more robust. The derivative is not likely to explode or vanish any more.

TABLE I
FULLY CONNECTED NETWORK TEST ACCURACY

| Optimizer | SGD | RMSProp | Adam |
|---|---|---|---|
| no batchnorm | 9.58% | 86.23% | 95.17% |
| batchnorm | 17.60% | 97.31% | 97.21% |

(3)Construct a Convolutional Neural Network and do the same things above.

Here we built a network with 7 convolution layers. Each layer's size is different from one another.

Again we compared the performance with and without batch normalization. Besides, we also tried different optimization methods, different conv layer numbers and different learning rate.

| Optimizer | SGD | RMSProp | Adam |
|---|---|---|---|
| no batchnorm | 60.65% | 98.85% | 98.98% |
| batchnorm | 80.69% | 99.02% | 98.82% |

| conv_layers | 3 | 5 | 7 |
|---|---|---|---|
| no batchnorm | 63.18% | 78.09% | 60.65% |
| batchnorm | 83.57% | 86.28% | 80.69% |

| learning_rate | 0.001 | 0.002 | 0.005 |
|---|---|---|---|
| no batchnorm | 66.94% | 79.41% | 93.04% |
| batchnorm | 81.25% | 88.98% | 94.79% |
| learning_rate | 0.01 | 0.02 | 0.1 |
| no batchnorm | 94.54% | 96.73% | 98.52% |
| batchnorm | 96.11% | 97.66% | 98.99% |

Easy to tell that batch normalization helps in most of the cases. And the reasons just like we said in (1)(2). Also, we may want to explain the batch normalization's role in CNN in another way. It is clear that batch normalization will make full use of more convolution kernel and layer, preventing some of them taking the dominant place and neglecting all the others.

Also for different optimization, SGD is the worst and RMSProp and Adam are really good.

For different structure of this network, the training results may be different. Intuitively, model with more layers has better approximation ability, but there are harder to train if play time is limited. Here we actually stopped before we can reach the best result since the loss is continuously decreasing, but we can also tell that some specific structure is better for this not-so-complex dataset.

Finally for the learning rate part, we can see at least for this model, we should use a larger learning rate. Also larger learning rate can prevent from being trapped in some local optimal but global sub-optimal.

Finally we can tell batch normalization helps a lot. Also RMSProp and Adam optimizer are good optimizer compared with SGD.
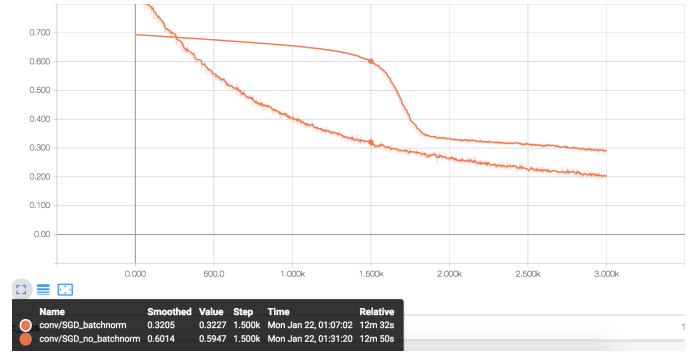


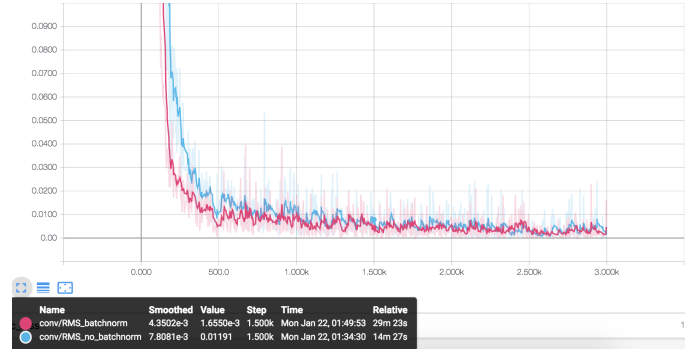Fig. 4. loss of CNN using SGD with or without batch norm



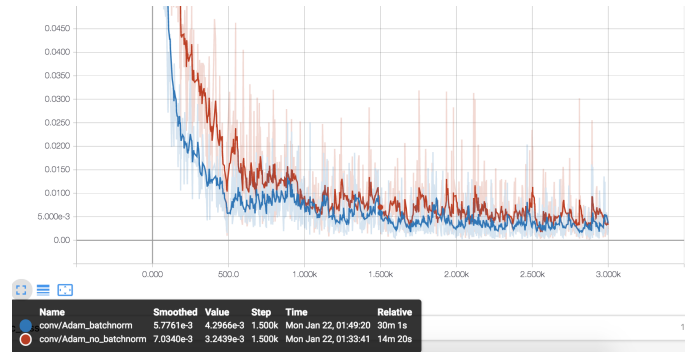Fig. 5. loss of CNN using RMSProp with or without batch norm



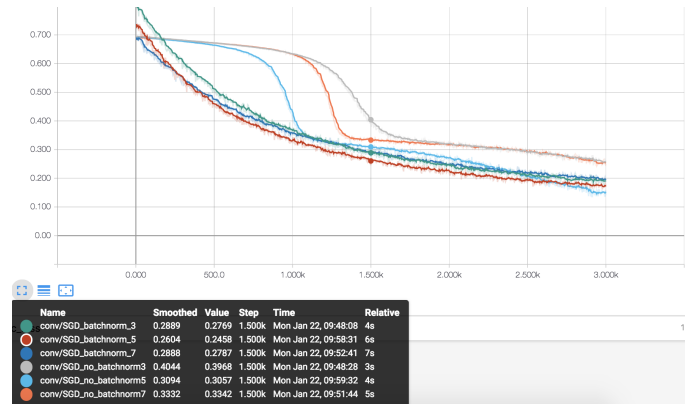Fig. 6. loss of CNN using Adam with or without batch norm



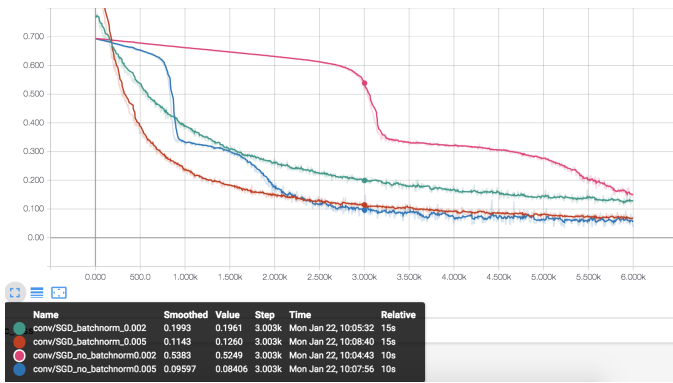Fig. 7. loss of CNN with different layers with or without batch norm

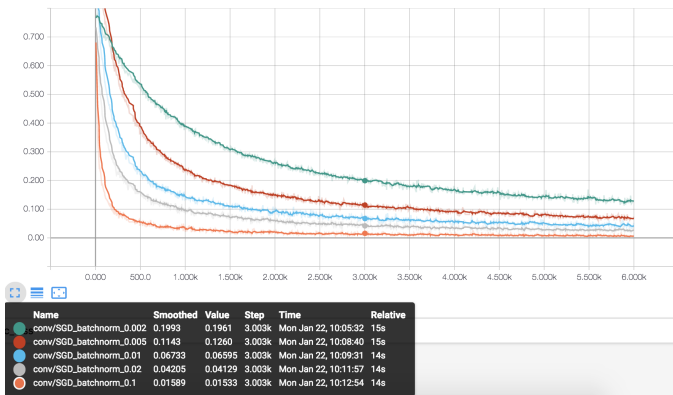Fig. 8. loss of CNN with different learning rate with or without batch norm



Fig. 9. loss of CNN with different learning rate with batch norm

# REFERENCES

[1] https://deepnotes.io/softmax-crossentropy
[2] Yes you should understand backprop - Andrej Karpathy - Medium
[3] https://en.wikipedia.org/wiki/Matrix_calculus
[4] https://deepnotes.io/batchnorm
[5] Ioffe, Sergey, and Christian Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift." International Conference on Machine Learning. 2015.
[6] TensorBoard: Visualizing Learning TensorFlow
[7] https://github.com/udacity/deep-learning/tree/master/batch-norm