hw 4

1.

$$x \xrightarrow{V} z_1 \xrightarrow{\text{ReLU (activate)}} a_1 \xrightarrow{W} z_2 \xrightarrow{\text{SoftMax (activate)}} a_2$$

input
hidden layer
+ bias

$(6000 \times 784) + 1)$   ~~785×200~~

Stochastic: $\in 785 \times 1$   $200 \times 1$   $201 \times 1$   $10 \times 1$   $10 \times 1$
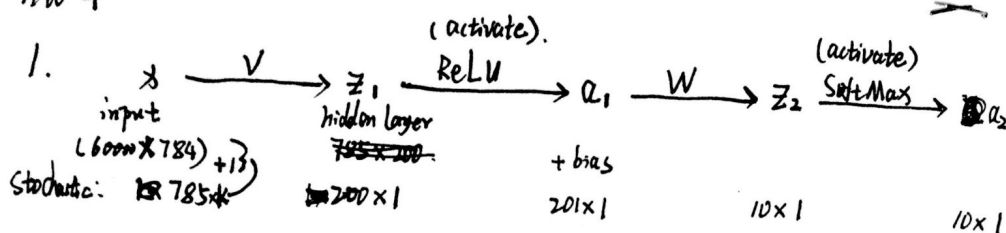
Loss: $J = -\sum_{k=1}^{10} y_k \ln z_k(x)$.

Softmax: $a_{2j} = g(z_{2j}) = \dfrac{e^{z_{2j}}}{\sum_{k=1}^{10} e^{z_{2k}}}$   $\dfrac{\partial a_{2j}}{\partial z_{2j}} = \begin{cases} a_{2j}(1 - a_{2j}), & i = j \\ -a_{2j} a_{2i}, & i \neq j \end{cases}$

use this result:

$$\Rightarrow \quad \frac{\partial L}{\partial z_{2i}} = -\left( y_i (1 - a_{2i}) + \sum_{k \neq i} -a_{2k} y_k \right)$$

$$= a_{2i} y_i - y_i + \sum_{k \neq i} a_{2k} y_k$$

$$= a_{2i}\left( \sum_i y_i \right) - y_i = a_{2i} - y_i$$

$$\Rightarrow \delta_i^L = \sum_i a_{2i} - y_i \quad \Rightarrow \underline{\delta^L = a_2 - y} \quad (10 \text{ by } 1).$$

$$\delta_i^{l-1} = \delta_i^L \cdot \frac{\partial z_{2i}}{\partial a_{1j}} \cdot \frac{\partial a_{1j}}{\partial z_{1k}} = (a_{2i} - y_i) \cdot W_{ij} \cdot \{0, 1\} \quad \text{if } a_{1j} > 0 \to 1$$

else $\to$ 1
if $a_{1j} \leq 0 \to 0$

$$\Rightarrow \delta^{l-1} = W^T \cdot \delta^L \cdot (\text{ReLU}'(a_1))^T \Rightarrow (\delta^L)^T \cdot W \cdot \text{ReLU}'(a_1)$$

$201 \times 10$   $10 \times 1$   $201 \times 200$   $1 \times 10$   $10 \times 201$   $201 \times 200$   $\Rightarrow 1 \times 200$

$$\Rightarrow \frac{\partial L}{\partial W} = \delta^L \cdot a_1^T$$

$10 \times 1$   $1 \times 201$

$$\Rightarrow \frac{\partial L}{\partial V} = (\delta^{l-1})^T \cdot x_i^T$$

$200 \times 1$   $1 \times 785$   $\Rightarrow 200 \times 785$

3. ① learning rate → 0.0001    minibatch size = 50    at beginning.
   decay rate → 0.95 every 1000 iteration.

   initialize W, V with. multivariate. normal. ( zero mean. $\delta = 0.05$ )
   use Nesterov accelerated gradient to increase speed of convergence.
   ( apply the momentum before calculating gradient)

   ②. final : 99.7% training accuracy. 99.6% validation accuracy.

   ③ Runing time :   first, I use batch=50, iter=10,000 → 20 mins.
                     Then increase batch size. $50 \to 100 \to 200 \to 500 \to 1000$.
                                              $\to 2000 \to 5000 \to 10000 \to 20000$
              for each batch size. I iterate 5 epoches.
              Although, I don't know why, the accuracy ~~was~~ increased after
              I added batch size. ( for batch =50, 97% is the highest,
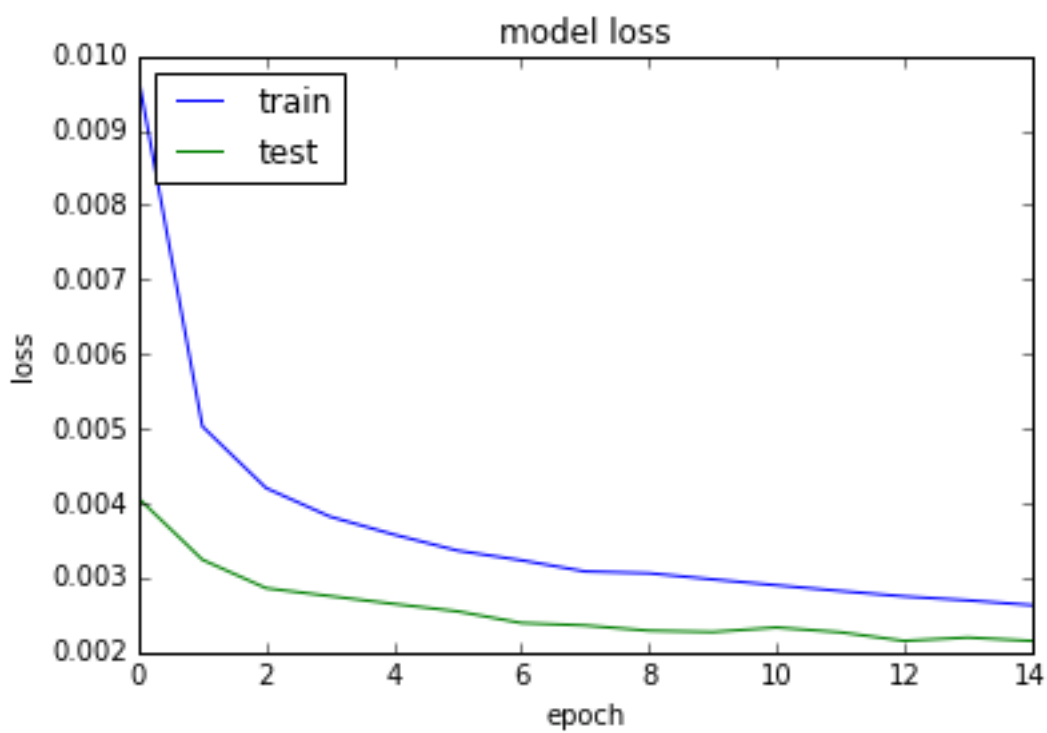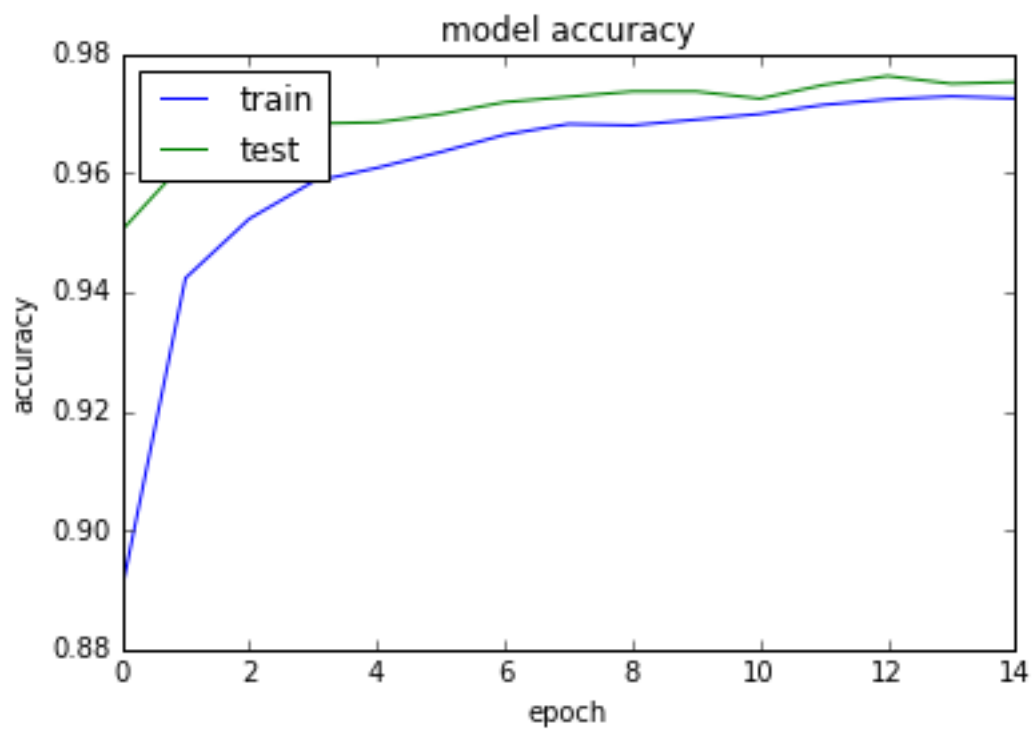                                        can't be improved any more).

   ④ plot and save.

4.    I used (1) L2 regularization. (2) use $\log(x+0.1)$ transformation to preprocess
      (3) Nesterov momentum. (4) to keep numerical stability. I used trick mentioned.
                in .pdf.    $\dfrac{e^{z_i} - max\{e^z\}}{\sum_{k}^{n} e^{z_k} - max\{e^z\}}$.

model accuracy

model loss

## Your Submissions

You are submitting as part of team ZhipengYu.    **Make a submission »**

**Note:** You can select up to **2** submissions to be used to calculate your final leaderboard score. If 2 submissions are not selected, they will be chosen based on your best submission scores on the public leaderboard.

Your final score will not be based on the same exact subset data as the public leaderboard, but rather a different private data subset of your full submission—your public score is only a rough indication of what your final score is. You should thus choose submissions that will most likely be best overall, and not necessarily just on the public subset.

**Your team's final score will be the best private submission score from the 2 selected submissions.**

| Submission | Files | Public Score | Selected? |
|---|---|---|---|
| Wed, 02 Nov 2016 10:13:11<br>Edit description | kaggle | 0.98480 | ☑ |
| Wed, 02 Nov 2016 09:20:23<br>Edit description | kaggle | 0.97580 | ☑ |