

Report of Homework3 IFT6390

Team Member

Qiang Ye (20139927), Lifeng Wan (20108546)

Coding Environment

python 3.5.2, numpy 1.14.2 matplotlib 2.2.0

1 THEORETICAL PART a (25 pts): derivatives and relationships between basic functions

Given

-- << logistic sigmoid >> $\text{sigmoid}(x) = \frac{1}{1+\exp(-x)}$

-- << hyperbolic tangend >> $\tanh(x) = \frac{\exp(x)-\exp(-x)}{\exp(x)+\exp(-x)}$

-- << soft plus >> $\text{softplus}(x) = \ln(1 + \exp(x))$

-- << sign >> function sign which returns +1 if its arguments is positive, -1 if negative and 0 if 0.

-- $\mathbf{1}_S(x)$ is the indicator function which returns 1 if $x \in S$ or x respects condition S , otherwise returns 0

-- << rectifier >> function which keeps only the positive part of its argument: $\text{rect}(x)$ returns x if $x \geq 0$ and returns 0 if $x < 0$. It is also named RELU (rectified linear unit): $\text{rect}(x) = \text{RELU}(x) = [x]_+ = \max(0, x) = x \cdot \mathbf{1}_{\{x>0\}}(x)$

1. Show that $\text{sigmoid}(x) = \frac{1}{2}(\tanh(\frac{1}{2}x) + 1)$

Answer

$$\begin{aligned}
 \frac{1}{2} \left(\tanh\left(\frac{1}{2}x\right) + 1 \right) &= \frac{1}{2} \left(\frac{\exp(\frac{1}{2}x) - \exp(-\frac{1}{2}x)}{\exp(\frac{1}{2}x) + \exp(-\frac{1}{2}x)} + 1 \right) \\
 &= \frac{1}{2} \left(\frac{2 \exp(\frac{1}{2}x)}{\exp(\frac{1}{2}x) + \exp(-\frac{1}{2}x)} \right) \\
 &= \frac{\exp(\frac{1}{2}x)}{\exp(\frac{1}{2}x) + \exp(-\frac{1}{2}x)} \\
 &= \frac{1}{1 + \frac{\exp(-\frac{1}{2}x)}{\exp(\frac{1}{2}x)}} \\
 &= \frac{1}{1 + \exp(-\frac{1}{2}x - \frac{1}{2}x)} \\
 &= \frac{1}{1 + \exp(-x)} \\
 &= \text{sigmoid}(x)
 \end{aligned}$$

2. Show that $\ln \text{sigmoid}(x) = -\text{softplus}(-x)$

Answer

$$\begin{aligned}\ln \text{sigmoid}(x) &= \ln \frac{1}{1 + \exp(-x)} \\ &= -\ln(1 + \exp(-x)) \\ &= -\text{softplus}(-x)\end{aligned}$$

3. Show that the derivative of the sigmoid is: $\text{sigmoid}'(x) = \frac{d \text{sigmoid}}{dx}(x) = \text{sigmoid}(x)(1 - \text{sigmoid}(x))$

Answer

$$\begin{aligned}\text{sigmoid}'(x) &= \frac{d \text{sigmoid}}{dx}(x) \\ &= \frac{d \frac{1}{1 + \exp(-x)}}{dx} \\ &= -\left(\frac{1}{1 + \exp(-x)}\right)^2 \frac{d(1 + \exp(-x))}{dx} \\ &= \frac{1}{(1 + \exp(-x))^2} \exp(-x) \\ &= \frac{1}{1 + \exp(-x)} \frac{\exp(-x)}{1 + \exp(-x)} \\ &= \frac{1}{1 + \exp(-x)} \frac{1 + \exp(-x) - 1}{1 + \exp(-x)} \\ &= \frac{1}{1 + \exp(-x)} \left(1 - \frac{1}{1 + \exp(-x)}\right) \\ &= \text{sigmoid}(x)(1 - \text{sigmoid}(x))\end{aligned}$$

4. Show that the tanh derivative is: $\tanh'(x) = 1 - \tanh^2(x)$

Answer

$$\begin{aligned}\tanh'(x) &= \frac{d \left(\frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)} \right)}{dx} \\ &= \frac{(\exp(x) + \exp(-x))(\exp(x) - \exp(-x))' - (\exp(x) - \exp(-x))(\exp(x) + \exp(-x))'}{(\exp(x) + \exp(-x))^2} \\ &= \frac{(\exp(x) + \exp(-x))^2 - (\exp(x) - \exp(-x))^2}{(\exp(x) + \exp(-x))^2} \\ &= 1 - \frac{(\exp(x) - \exp(-x))^2}{(\exp(x) + \exp(-x))^2} \\ &= 1 - \tanh^2(x)\end{aligned}$$

5. Write the sign function using only indicator functions: $\text{sign}(x) = \dots$

Answer $\text{sign}(x) = \mathbf{1}_{\{x>0\}}(x) - \mathbf{1}_{\{x<0\}}(x)$

6. Write the derivative of the absolute function $\text{abs}(x) = |x|$. Note: its derivative at 0 is not defined, but your function abs' can return 0 at 0. Note2: use the sign function: $\text{abs}'(x) = \dots$

Answer $\text{abs}'(x) = \text{sign}(x)$

7. Write the derivative of the function rect . Note: its derivative at 0 is undefined, but your function can return 0 at 0. Note2: use the indicator function. $\text{rect}'(x) = \dots$

Answer $\text{rect}'(x) = \mathbf{1}_{\{x>0\}}(x)$

8. Let the squared L_2 norm of a vector be: $\|\mathbf{x}\|_2^2 = \sum_i \mathbf{x}_i^2$. Write the vector of the gradient: $\frac{\partial \|\mathbf{x}\|_2^2}{\partial \mathbf{x}} = \dots$

Answer

$$\frac{\partial \|\mathbf{x}\|_2^2}{\partial \mathbf{x}} = \left(\frac{\partial \sum_i \mathbf{x}_i^2}{\partial \mathbf{x}_1}, \frac{\partial \sum_i \mathbf{x}_i^2}{\partial \mathbf{x}_2}, \dots \right)^T = (2\mathbf{x}_1, 2\mathbf{x}_2, \dots)^T = 2\mathbf{x}$$

9. Let the norm L_1 of a vector be: $\|\mathbf{x}\|_1 = \sum_i |\mathbf{x}_i|$. Write the vector of the gradient: $\frac{\partial \|\mathbf{x}\|_1}{\partial \mathbf{x}} = \dots$

Answer

$$\begin{aligned} \frac{\partial \|\mathbf{x}\|_1}{\partial \mathbf{x}} &= \left(\frac{\partial \sum_i |\mathbf{x}_i|}{\partial \mathbf{x}_1}, \frac{\partial \sum_i |\mathbf{x}_i|}{\partial \mathbf{x}_2}, \dots \right)^T = (\text{sign}(\mathbf{x}_1), \text{sign}(\mathbf{x}_2), \dots)^T \\ &= \text{sign}(\mathbf{x}) \end{aligned}$$

2 THEORETICAL PART b (25 pts): Gradient computation for parameters optimization in a neural net for multiclass classification

Let $D_n = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$ be the dataset with $\mathbf{x}^{(i)} \in \mathbb{R}^d$ and $y^{(i)} \in \{1, \dots, m\}$ indicating the class within m classes. **For vectors and matrices in the following equations, vectors are by default considered to be column vectors.**

Consider a neural net of the type *Multilayer perceptron* (MLP) with only one hidden layer (meaning 3 layers total if we count the input and output layers). The hidden layer is made of d_h neurons fully connected to the input layer. We shall consider a non linearity of type **rectifier** (Rectified Linear Unit of **RELU**) for the hidden layer. The output layer is made of m neurons that are fully connected to the hidden layer. They are equipped with a **softmax** non linearity. The output of the j^{th} neuron of the output layer gives a score for the class j which is interpreted as the probability of \mathbf{x} being of class j .

It is highly recommended that you draw the neural net as it helps understanding all the steps.

1. Let $\mathbf{W}^{(1)}$ a $d_h \times d$ matrix of weights and $\mathbf{b}^{(1)}$ the bias vector be the connections between the input layer and the hidden layer. What is the dimension of $\mathbf{b}^{(1)}$? Give the formula of the preactivation vector (before the non linearity) of the neurons of hidden layer \mathbf{h}^a given \mathbf{x} as input, first in a matrix form ($\mathbf{h}^a = \dots$), and then details on how to compute one element $\mathbf{h}_j^a = \dots$. Write the output vector of the hidden layer \mathbf{h}^s with respect to \mathbf{h}^a .

Answer

$$\begin{aligned} \mathbf{b}^{(1)} &\in \mathbb{R}^{d_h} \\ \mathbf{h}^a &= \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)} \\ \mathbf{h}_j^a &= \mathbf{W}_j^{(1)}\mathbf{x} + \mathbf{b}_j^{(1)} = \sum_{i=1}^d \mathbf{W}_{ji}^{(1)}\mathbf{x}_i + \mathbf{b}_j^{(1)} \\ \mathbf{h}^s &= \text{RELU}(\mathbf{h}^a) \end{aligned}$$

2. Let $\mathbf{W}^{(2)}$ a weight matrix and $\mathbf{b}^{(2)}$ a bias vector be the connections between the hidden layer and the output layer. What are the dimensions of $\mathbf{W}^{(2)}$ and $\mathbf{b}^{(2)}$? Give the formula of the activation function of the neurons of the output layer \mathbf{o}^a with respect to input \mathbf{h}^s in a matrix form and then write in a detailed form for \mathbf{o}_k^a .

Answer

$$\begin{aligned} \mathbf{W}^{(2)} &\text{ is a matrix of } m \times d_h, \mathbf{b}^{(2)} \in \mathbb{R}^m. \\ \mathbf{o}^a &= \mathbf{W}^{(2)}\mathbf{h}^s + \mathbf{b}^{(2)} \\ \mathbf{o}_k^a &= \mathbf{W}_k^{(2)}\mathbf{h}^s + \mathbf{b}_k^{(2)} = \sum_{j=1}^{d_h} \mathbf{W}_{kj}^{(2)}\mathbf{h}_j^s + \mathbf{b}_k^{(2)} \end{aligned}$$

3. The output of the neurons at the output layer is given by:

$$\mathbf{o}^s = \text{softmax}(\mathbf{o}^a)$$

Give the precise equation for \mathbf{o}_k^s using the softmax (formular with the exp). **Show** that the \mathbf{o}_k^s are positive and sum to 1. Why is this important?

Answer

$$\mathbf{o}_k^s = \frac{\exp(\mathbf{o}_k^a)}{\sum_{i=1}^m \exp(\mathbf{o}_i^a)}$$

\mathbf{o}_k^s is always positive because $\exp(\mathbf{o}_k^a)$ and $\sum_k \exp(\mathbf{o}_k^a)$ are always positive with $\mathbf{o}_k^a \in \mathbb{R}$.

$$\sum_{k=1}^m \mathbf{o}_k^s = \sum_{k=1}^m \frac{\exp(\mathbf{o}_k^a)}{\sum_{i=1}^m \exp(\mathbf{o}_i^a)} = \frac{\sum_{k=1}^m \exp(\mathbf{o}_k^a)}{\sum_{i=1}^m \exp(\mathbf{o}_i^a)} = 1$$

softmax function maps a score vector to a vector in which all elements sum to 1 such that each value in the vector can represent the probability with which an example belongs to a certain class. Therefore, it is crucial that all probabilities sum to 1 since an example should always be in one of the m classes.

4. The neural net computes, for an input vector \mathbf{x} , a vector of probability scores $\mathbf{o}^s(\mathbf{x})$. The probability, computed by a neural net, that an observation \mathbf{x} belongs to class y is given by the y^{th} output $\mathbf{o}_y^s(\mathbf{x})$. This suggests a loss function such as:

$$L(\mathbf{x}, y) = -\log \mathbf{o}_y^s(\mathbf{x})$$

Find the equation of L as a function of the vector \mathbf{o}^a . It is easily achievable with the correct substitution using the equation of the previous question.

Answer

$$\begin{aligned} L(\mathbf{x}, y) &= -\log \mathbf{o}_y^s(\mathbf{x}) = -\log \frac{\exp(\mathbf{o}_y^a(\mathbf{x}))}{\sum_{k=1}^m \exp(\mathbf{o}_k^a(\mathbf{x}))} \\ &= \log \sum_{k=1}^m \exp(\mathbf{o}_k^a(\mathbf{x})) - \log \exp(\mathbf{o}_y^a(\mathbf{x})) \\ &= \log \sum_{k=1}^m \exp(\mathbf{o}_k^a(\mathbf{x})) - \mathbf{o}_y^a(\mathbf{x}) \end{aligned}$$

5. The training of the neural net will consist of finding parameters that minimize the empirical risk \hat{R} associated with this loss function. What is \hat{R} ? What is precisely the set θ of parameters of the network? How many scalar parameters n_θ are there? Write down the optimization problem of training the network in order to find the optimal values for these parameters.

Answer

$$\hat{R} = \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}^{(i)}, y^{(i)})$$

$$\theta = \{(\mathbf{W}^{(1)}, \mathbf{b}^{(1)}), (\mathbf{W}^{(2)}, \mathbf{b}^{(2)})\}$$

$$n_\theta = d_h \times (d + 1) + m \times (d_h + 1)$$

The optimization process is to find an optimal θ^* :

$$\theta^* = \arg \min_{\theta} \hat{R} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}^{(i)}, y^{(i)})$$

6. To find a solution to this optimization problem, we will use gradient descent. What is the (batch) gradient descent equation for this problem?

Answer

$$\theta = \theta - \alpha \frac{\partial \hat{R}(D_n)}{\partial \theta}$$

where α is a learning rate (step size) when performing gradient descent.

7. We can compute the vector of the gradient of the empirical risk \hat{R} with respect to the parameters set θ this way

$$\begin{pmatrix} \frac{\partial \hat{R}}{\partial \theta_1} \\ \vdots \\ \frac{\partial \hat{R}}{\partial \theta_{n\theta}} \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \frac{\partial L(\mathbf{x}_i, y_i)}{\partial \theta_1} \\ \vdots \\ \frac{\partial L(\mathbf{x}_i, y_i)}{\partial \theta_{n\theta}} \end{pmatrix}$$

This hints that we only need to know how to compute the gradient of the loss L with an example (\mathbf{x}, y) with respect to the parameters, defined as followed:

$$\frac{\partial L}{\partial \theta} = \begin{pmatrix} \frac{\partial L}{\partial \theta_1} \\ \vdots \\ \frac{\partial L}{\partial \theta_{n\theta}} \end{pmatrix} = \begin{pmatrix} \frac{\partial L(\mathbf{x}, y)}{\partial \theta_1} \\ \vdots \\ \frac{\partial L(\mathbf{x}, y)}{\partial \theta_{n\theta}} \end{pmatrix}$$

We shall use **gradient backpropagation**, starting with loss L and going to the output layer \mathbf{o} then down the hidden layer \mathbf{h} then finally at the input layer \mathbf{x} .

Show that

$$\frac{\partial L}{\partial \mathbf{o}^a} = \mathbf{o}^s - \text{onehot}_m(y)$$

Note: Start from the expression of L as a function of \mathbf{o}^a that you previously found. Start by computing $\frac{\partial L}{\partial \mathbf{o}_k^a}$ for $k \neq y$ (using the start of the expression of the logarithm derivate). Do the same thing for $\frac{\partial L}{\partial \mathbf{o}_y^a}$.

Answer

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{o}^a} &= \frac{\partial}{\partial \mathbf{o}^a} \left(\log \sum_{k=1}^m \exp(\mathbf{o}_k^a) - \mathbf{o}_y^a \right) \\ &= \begin{pmatrix} \frac{\partial}{\partial \mathbf{o}_1^a} \left(\log \sum_{k=1}^m \exp(\mathbf{o}_k^a) - \mathbf{o}_y^a \right) \\ \vdots \\ \frac{\partial}{\partial \mathbf{o}_m^a} \left(\log \sum_{k=1}^m \exp(\mathbf{o}_k^a) - \mathbf{o}_y^a \right) \end{pmatrix} \end{aligned}$$

when $k \neq y$,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{o}_k^a} \left(\log \sum_{k=1}^m \exp(\mathbf{o}_k^a) - \mathbf{o}_y^a \right) &= \frac{\partial}{\partial \mathbf{o}_k^a} \left(\log \sum_{k=1}^m \exp(\mathbf{o}_k^a) \right) \\ &= \frac{1}{\sum_{k=1}^m \exp(\mathbf{o}_k^a)} \frac{\partial}{\partial \mathbf{o}_k^a} \sum_{k=1}^m \exp(\mathbf{o}_k^a) \\ &= \frac{\exp(\mathbf{o}_k^a)}{\sum_{k=1}^m \exp(\mathbf{o}_k^a)} \\ &= \mathbf{o}_k^s \end{aligned}$$

when $k = y$,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{o}_k^a} \left(\log \sum_{k=1}^m \exp(\mathbf{o}_k^a) - \mathbf{o}_y^a \right) &= \frac{\partial}{\partial \mathbf{o}_k^a} \left(\log \sum_{k=1}^m \exp(\mathbf{o}_k^a) \right) - 1 \\ &= \frac{1}{\sum_{k=1}^m \exp(\mathbf{o}_k^a)} \left(\frac{\partial}{\partial \mathbf{o}_k^a} \sum_{k=1}^m \exp(\mathbf{o}_k^a) \right) - 1 \\ &= \frac{\exp(\mathbf{o}_k^a)}{\sum_{k=1}^m \exp(\mathbf{o}_k^a)} - 1 \\ &= \mathbf{o}_k^s - 1 \end{aligned}$$

Combine two above equations, we can conclude that:

$$\frac{\partial L}{\partial \mathbf{o}^a} = \mathbf{o}^s - \text{onehot}_m(y)$$

8. What is the numpy equivalent expression (it can fit 2 operations)?

`grad_oa = ...`

...

Answer

numpy form:

for one example

`grad_oa = np.copy(os) # (n_output, 1)`

`np.put(grad_oa, y, grad_oa[y,0]-1) # (n_output, 1)`

for batch_size

`def one_hot(m, y):`

`y = y.reshape(-1,) # (batch_size,)`

`batch_size = y.shape[0]`

`result = np.zeros((batch_size, m))`

`result[np.arange(batch_size), y] = 1`

`return result`

`grad_oa = o_s - one_hot(m, y) # (batch_size, n_output)`

IMPORTANT: From now on when we ask to "compute" the gradients or partial derivatives, you only need to write them as function of previously computed derivatives (**do not substitute the whole expressions already computed in the previous questions!**)

9. Compute the gradients with respect to parameters $\mathbf{W}^{(2)}$ and $\mathbf{b}^{(2)}$ of the output layer. Since L depends on $\mathbf{W}_{kj}^{(2)}$ and $\mathbf{b}_k^{(2)}$ only through \mathbf{o}_k^a the result of the chain rule is:

$$\frac{\partial L}{\partial \mathbf{W}_{kj}^{(2)}} = \frac{\partial L}{\partial \mathbf{o}_k^a} \frac{\partial \mathbf{o}_k^a}{\partial \mathbf{W}_{kj}^{(2)}}$$

and

$$\frac{\partial L}{\partial \mathbf{b}_k^{(2)}} = \frac{\partial L}{\partial \mathbf{o}_k^a} \frac{\partial \mathbf{o}_k^a}{\partial \mathbf{b}_k^{(2)}}$$

Answer

$$\frac{\partial L}{\partial \mathbf{W}_{kj}^{(2)}} = \frac{\partial L}{\partial \mathbf{o}_k^a} \frac{\partial}{\partial \mathbf{W}_{kj}^{(2)}} \left(\sum_{j=1}^{d_h} \mathbf{W}_{kj}^{(2)} \mathbf{h}_j^s + \mathbf{b}_k^{(2)} \right) = \frac{\partial L}{\partial \mathbf{o}_k^a} \mathbf{h}_j^s$$

$$\frac{\partial L}{\partial \mathbf{b}_k^{(2)}} = \frac{\partial L}{\partial \mathbf{o}_k^a} \frac{\partial}{\partial \mathbf{b}_k^{(2)}} \left(\sum_{j=1}^{d_h} \mathbf{W}_{kj}^{(2)} \mathbf{h}_j^s + \mathbf{b}_k^{(2)} \right) = \frac{\partial L}{\partial \mathbf{o}_k^a}$$

10. Write down the gradient of the last question in matrix form and define the dimensions of all matrix or vectors involved.

(What are the dimensions?)

Take time to understand why the above equalities are the same as the equations of the last question.

Give the numpy form:

grad_b2 = ...

grad_W2 = ...

Answer

$$\frac{\partial L}{\partial \mathbf{W}^{(2)}} = \frac{\partial L}{\partial \mathbf{o}^a} (\mathbf{h}^s)^T$$

$$\frac{\partial L}{\partial \mathbf{b}^{(2)}} = \frac{\partial L}{\partial \mathbf{o}^a}$$

where,

 \mathbf{o}^a is a vector of $(m \times 1)$ \mathbf{h}^s is a vector of $(d_h \times 1)$ $\mathbf{W}^{(2)}$ is a matrix of $(m \times d_h)$ $\mathbf{b}^{(2)}$ is a vector of $(m \times 1)$ $\frac{\partial L}{\partial \mathbf{W}^{(2)}}$ is a matrix of $(m \times d_h)$ $\frac{\partial L}{\partial \mathbf{b}^{(2)}}$ is a vector of $(m \times 1)$ $\frac{\partial L}{\partial \mathbf{o}^a}$ is a vector of $(m \times 1)$

numpy form:

for one example

grad_W2 = np.dot(grad_oa, hs.T) # (n_output, n_hidden)

grad_b2 = grad_oa # (n_output, 1)

for batch_size

grad_W2 = np.dot(grad_oa.T, hs) # (n_output, n_hidden)

grad_b2 = grad_oa # (batch_size, n_output)

grad_b2 = np.mean(grad_b2, axis = 0, keepdims = True).T # (n_output, 1)

11. What is the partial derivate of the loss L with respect to the output of the neurons at the hidden layer?Since L depends on \mathbf{h}_j^s only through the activations of the output neurons \mathbf{o}^a the chain rule yields:

$$\frac{\partial L}{\partial \mathbf{h}_j^s} = \sum_{k=1}^m \frac{\partial L}{\partial \mathbf{o}_k^a} \frac{\partial \mathbf{o}_k^a}{\partial \mathbf{h}_j^s}$$

Answer

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{h}_j^s} &= \sum_{k=1}^m \frac{\partial L}{\partial \mathbf{o}_k^a} \frac{\partial \mathbf{o}_k^a}{\partial \mathbf{h}_j^s} \\ &= \sum_{k=1}^m \frac{\partial L}{\partial \mathbf{o}_k^a} \frac{\partial}{\partial \mathbf{h}_j^s} \left(\sum_{j=1}^{d_h} \mathbf{W}_{kj}^{(2)} \mathbf{h}_j^s + \mathbf{b}_k^{(2)} \right) \\ &= \sum_{k=1}^m \frac{\partial L}{\partial \mathbf{o}_k^a} \mathbf{W}_{kj}^{(2)} \end{aligned}$$

12. Write down the gradient of the last question in matrix form and define the dimensions of all matrix or vectors involved.

(What are the dimensions?)

Take time to understand why the above equalities are the same as the equations of the last question.

Give the numpy form: `grad_hs = ...`

Answer

$$\frac{\partial L}{\partial \mathbf{h}^s} = (\mathbf{W}^{(2)})^T \frac{\partial L}{\partial \mathbf{o}^a}$$

where,

\mathbf{h}^s is a vector of $d_h \times 1$

$\mathbf{W}^{(2)}$ is a matrix of $m \times d_h$

$\frac{\partial L}{\partial \mathbf{o}^a}$ is a vector of $m \times 1$

$\frac{\partial L}{\partial \mathbf{h}^s}$ is a vector of $d_h \times 1$

numpy form:

```
# for one example
grad_hs = np.dot(W2.T, grad_oa)

# for batch_size
grad_hs = np.dot(grad_oa, self.W2) # (batch_size, n_hidden)
```

13. What is the partial derivate of the loss L with respect to the activation of the neurons at the hidden layer?

Since L depends on the activation \mathbf{h}_j^a only through \mathbf{h}_j^s of this neuron, the chain rule gives:

$$\frac{\partial L}{\partial \mathbf{h}_j^a} = \frac{\partial L}{\partial \mathbf{h}_j^s} \frac{\partial \mathbf{h}_j^s}{\partial \mathbf{h}_j^a}$$

Note $\mathbf{h}_j^s = \text{rect}(\mathbf{h}_j^a)$: the rectifier function is applied elementwise. Start by writing the derivate of the rectifier function $\frac{\partial \text{rect}(z)}{\partial z} = \text{rect}'(z) = \dots$

Answer

$$\frac{\partial L}{\partial \mathbf{h}_j^a} = \frac{\partial L}{\partial \mathbf{h}_j^s} \frac{\partial \mathbf{h}_j^s}{\partial \mathbf{h}_j^a} = \frac{\partial L}{\partial \mathbf{h}_j^s} \mathbf{1}_{\{\mathbf{h}_j^a > 0\}}(\mathbf{h}_j^a)$$

14. Write down the gradient of the last question in matrix form and define the dimensions of all matrix or vectors involved. Give the numpy form.

Answer

$$\frac{\partial L}{\partial \mathbf{h}^a} = \frac{\partial L}{\partial \mathbf{h}^s} \odot \mathbf{1}_{\{\mathbf{h}^a > 0\}}(\mathbf{h}^a)$$

where \odot denotes the elementwise multiplication. Each variable in the above equation is a vector of $d_h \times 1$.

numpy form:

```
# for one example
grad_ha = np.multiply(grad_hs, np.maximum(np.sign(ha),0)) # (n_hidden, 1)

# for batch_size:
grad_ha = np.multiply(grad_hs, np.maximum(np.sign(ha),0)) # (batch_size, n_hidden)
```


15. What is the gradient with respect to the parameters $\mathbf{W}^{(1)}$ and $\mathbf{b}^{(1)}$ of the hidden layer?

Note: same logic as a previous question

Answer

$$\frac{\partial L}{\partial \mathbf{W}_{ji}^{(1)}} = \frac{\partial L}{\partial \mathbf{h}_j^a} \frac{\partial \mathbf{h}_j^a}{\partial \mathbf{W}_{ji}^{(1)}} = \frac{\partial L}{\partial \mathbf{h}_j^a} \frac{\partial}{\partial \mathbf{W}_{ji}^{(1)}} \left(\sum_{i=1}^d \mathbf{W}_{ji}^{(1)} \mathbf{x}_i + \mathbf{b}_j^{(1)} \right) = \mathbf{x}_i \frac{\partial L}{\partial \mathbf{h}_j^a}$$

$$\frac{\partial L}{\partial \mathbf{b}_j^{(1)}} = \frac{\partial L}{\partial \mathbf{h}_j^a} \frac{\partial}{\partial \mathbf{b}_j^{(1)}} \left(\sum_{i=1}^d \mathbf{W}_{ji}^{(1)} \mathbf{x}_i + \mathbf{b}_j^{(1)} \right) = \frac{\partial L}{\partial \mathbf{h}_j^a}$$

16. Write down the gradient of the last question in matrix form and define the dimensions of all matrix or vectors involved. Give the numpy form.

Note: same logic as a previous question

Answer

$$\frac{\partial L}{\partial \mathbf{W}^{(1)}} = \frac{\partial L}{\partial \mathbf{h}^a} \mathbf{x}^T$$

$$\frac{\partial L}{\partial \mathbf{b}^{(1)}} = \frac{\partial L}{\partial \mathbf{h}^a}$$

where,

\mathbf{h}^a is a vector of $(d_h \times 1)$

\mathbf{x} is a vector of $(d \times 1)$

$\mathbf{W}^{(1)}$ is a matrix of $(d_h \times d)$

$\mathbf{b}^{(1)}$ is a vector of $(d_h \times 1)$

$\frac{\partial L}{\partial \mathbf{W}^{(1)}}$ is a matrix of $(d_h \times d)$

$\frac{\partial L}{\partial \mathbf{b}^{(1)}}$ is a vector of $(d_h \times 1)$

$\frac{\partial L}{\partial \mathbf{h}^a}$ is a vector of $(d_h \times 1)$

numpy form:

```
# for one example
grad_W1 = np.dot(grad_ha, x.T) # (n_hidden, n_input)
grad_b1 = grad_ha # (n_hidden, 1)

# for batch_size
grad_W1 = np.dot(grad_ha.T, X) # (n_hidden, n_input)
grad_b1 = grad_ha # (batch_size, n_hidden)
grad_b1 = np.mean(grad_b1, axis = 0, keepdims = True).T # (n_hidden, 1)
```

17. What are the partial derivatives of the loss L with respect to \mathbf{x} ?

Note: same logic as a previous question

Answer

$$\begin{aligned}
\frac{\partial L}{\partial \mathbf{x}_i} &= \sum_{j=1}^{d_h} \frac{\partial L}{\partial \mathbf{h}_j^a} \frac{\partial \mathbf{h}_j^a}{\partial \mathbf{x}_i} \\
&= \sum_{j=1}^{d_h} \frac{\partial L}{\partial \mathbf{h}_j^a} \frac{\partial}{\partial \mathbf{x}_i} \left(\sum_{i=1}^d \mathbf{W}_{ji}^{(1)} \mathbf{x}_i + \mathbf{b}_j^{(1)} \right) \\
&= \sum_{j=1}^{d_h} \frac{\partial L}{\partial \mathbf{h}_j^a} \mathbf{W}_{ji}^{(1)}
\end{aligned}$$

Matrix form:

$$\frac{\partial L}{\partial \mathbf{x}} = (\mathbf{W}^{(1)})^T \frac{\partial L}{\partial \mathbf{h}^a}$$

where,

\mathbf{x} is a vector of $d \times 1$

$\mathbf{W}^{(1)}$ is a matrix of $d_h \times d$

$\frac{\partial L}{\partial \mathbf{h}^a}$ is a vector of $d_h \times 1$

$\frac{\partial L}{\partial \mathbf{x}}$ is a vector of $d \times 1$

numpy form:

```
# for one example
grad_x = np.dot(W1.T, grad_ha) # (n_input, 1)

# for batch_size
grad_x = np.dot(grad_ha, self.W1) # (batch_size, n_input)
```

18. We will now consider a **regularized** empirical risk: $\tilde{R} = \hat{R} + \mathcal{L}(\theta)$, where θ is the vector of all the parameters in the network and $\mathcal{L}(\theta)$ describes a scalar penalty as a function of the parameters θ . The penalty is given importance according to a prior preferences for the values of θ . The L_2 (quadratic) regularization that penalizes the square norm (norm L_2) of the weights (but not the biases) is more standard, is used in ridge regression and is sometimes called "weight-decay". Here we shall consider a double regularization L_2 and L_1 which is sometimes named "elastic net" and we will use different **hyperparameters** (positive scalars $\lambda_{11}, \lambda_{12}, \lambda_{21}, \lambda_{22}$) to control the effect of the regularization at each layer

$$\begin{aligned}
\mathcal{L}(\theta) &= \mathcal{L}(\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}) \\
&= \lambda_{11} \|\mathbf{W}^{(1)}\|_1 + \lambda_{12} \|\mathbf{W}^{(1)}\|_2^2 + \lambda_{21} \|\mathbf{W}^{(2)}\|_1 + \lambda_{22} \|\mathbf{W}^{(2)}\|_2^2 \\
&= \lambda_{11} \left(\sum_{ij} |\mathbf{W}_{ij}^{(1)}| \right) + \lambda_{12} \left(\sum_{ij} (\mathbf{W}_{ij}^{(1)})^2 \right) + \lambda_{21} \left(\sum_{ij} |\mathbf{W}_{ij}^{(2)}| \right) + \lambda_{22} \left(\sum_{ij} (\mathbf{W}_{ij}^{(2)})^2 \right)
\end{aligned}$$

We will in fact minimize the regularized risk \tilde{R} instead of \hat{R} . How does this change the gradient with respect to the different parameters?

Answer

$$\begin{aligned}
\frac{\partial \tilde{R}}{\partial \mathbf{W}_{ij}^{(1)}} &= \frac{\partial \hat{R}}{\partial \mathbf{W}_{ij}^{(1)}} + \lambda_{11} \text{sign}(\mathbf{W}_{ij}^{(1)}) + 2\lambda_{12} \mathbf{W}_{ij}^{(1)} \\
\frac{\partial \tilde{R}}{\partial \mathbf{W}_{ij}^{(2)}} &= \frac{\partial \hat{R}}{\partial \mathbf{W}_{ij}^{(2)}} + \lambda_{21} \text{sign}(\mathbf{W}_{ij}^{(2)}) + 2\lambda_{22} \mathbf{W}_{ij}^{(2)} \\
\frac{\partial \tilde{R}}{\partial \mathbf{b}_i^{(2)}} &= \frac{\partial \hat{R}}{\partial \mathbf{b}_i^{(2)}}
\end{aligned}$$

$$\frac{\partial \tilde{R}}{\partial \mathbf{b}_i^{(1)}} = \frac{\partial \hat{R}}{\partial \mathbf{b}_i^{(1)}},$$

PRACTICAL PART (50 pts) : Neural net-work implementation and experiments

We ask you to implement a neural network where you compute the gradients using the formulas derived in the previous part (including elastic net type regularization). You must not use an existing neural network library, but you must use the derivation of part 2 (with corresponding variable names, etc). Note that you can reuse the general learning algorithm structure that we used in the demos, as well as the functions used to plot the decision functions.

Useful details on implementation : (cutted here. See homework description for detail)

Brief Description of the Answer / Code

We first went through all the experiments demanded. After knowing what we will challenge, we decided first to do some basic works that all experiments will need, including a complete implementation of our `MLP` class who performs some basic functions for the experiments. After that, for each specific experiment, we first gave a detail explanation about related member functions of `MLP` we used for that experiment, then we wrote extra codes based on the member functions, run it to show the results. So, the whole practical part is divided into two main parts: **Preparation of the experiments** and **Experiments and Results**.

Preparation of the experiments

1. Implementation of two basic functions

We carefully implemented two functions: `my_softmax(z)` and `one_hot(m, y)` to do the **softmax** prediction and **one_hot** encoding, both of which are implemented based on Matrix operation with one example also supported as an `batch_size` of 1.

In [1]:

```

1 %matplotlib inline
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 def my_softmax(z):
6     return np.exp(z) / np.sum(np.exp(z), axis=1, keepdims=True)
7
8 def one_hot(m, y):
9     y = y.reshape(-1, ) # (batch_size, )
10    y = y.astype(int)
11    batch_size = y.shape[0]
12    result = np.zeros((batch_size, m))
13    result[np.arange(batch_size), y] = 1
14    return result
15
16 DEBUG = True
17 def debug(*kargs):
18     global DEBUG
19     if DEBUG:
20         print(*kargs)

```

2. Implementation of MLP class

Note:

1. Initialization of the network parameters(W_1 , W_2 , b_1 , b_2) is not performed in the function `__init__` as we might not know the number of features of input and output until we know the X and y .
2. Two set of methods (loop and matrix expression) about the forward and backward propagation are implemented as required by the experiments which will be explained in detail in the answer of related experiments.

The completed implementation of class `MLP` is as follows:

In [2]:

```

1  class MLP():
2      """one hidden layer with RELU and Softmax
3      """
4      def __init__(self,
5                  n_hidden = 2,
6                  lambdas = [0, 0, 0, 0],
7                  learning_rate = 1e-3,
8                  epochs = 200,
9                  batch_size = 10,
10                 verbose = False,
11                 n_input = None,
12                 n_output = None,
13                 early_stop = False
14                 ):
15         self.n_i = n_input
16         self.n_o = n_output
17         self.n_h = n_hidden
18         self.lambdas = lambdas
19         self.lr = learning_rate
20         self.epochs = epochs
21         self.batch_size = batch_size
22         self.early_stop = early_stop
23         self.verbose = verbose
24
25
26     def _init_weights_biases(self, X, y, n_output = None):
27         """number of input and output neurons are decided when dataset
28         is known, initialization of parameters are performed here.
29         Params
30             X: input of training data ndarray (sample_size, n_features)
31             y: true class labels of X ndarray (sample_size, 1)
32             n_output: number of true classes, is None, can be calculated
33         from y.
34         Returns
35             None
36         """
37         if n_output is None:
38             y = y.reshape(-1, )
39             self.n_o = int(np.max(y) + 1)
40         else:
41             self.n_o = n_output
42         np.random.seed(0) # for compare
43         sample_size, n_input = X.shape
44         self.n_i = n_input
45         _high = 1 / np.sqrt(self.n_i)
46         self.W1 = np.random.uniform(-1 * _high, _high,
47                                     (self.n_h, self.n_i))
48         self.b1 = np.zeros((self.n_h, 1)) # column vector
49
50         _high = 1 / np.sqrt(self.n_h)
51         self.W2 = np.random.uniform(-1 * _high, _high,
52                                     (self.n_o, self.n_h))
53         self.b2 = np.zeros((self.n_o, 1)) # column vector
54
55
56     def _init_grads(self):
57         """Initialization of gradient of parameters
58         Params:None
59         Returns

```

```

60         grads: zero gradient of all parameters dict
61     """
62     grad_W1 = np.zeros_like(self.W1)
63     grad_W2 = np.zeros_like(self.W2)
64     grad_b1 = np.zeros_like(self.b1)
65     grad_b2 = np.zeros_like(self.b2)
66     grads = {"grad_W1": grad_W1, "grad_W2": grad_W2,
67             "grad_b1": grad_b1, "grad_b2": grad_b2,
68             }
69     return grads
70
71
72 def _fprop_one_example(self, x):
73     """forward propagation of the network using one example.
74     Params
75         x: example ndarray (n_features, 1)
76     Returns
77         cache: a dict variable including computation results used for
78         back propagation.
79     """
80     x = x.reshape(-1, 1) # column vector (n_input, 1)
81     ha = np.dot(self.W1, x) + self.b1 # (n_hidden, 1)
82     hs = np.maximum(ha, 0) # (n_hidden, 1)
83
84     oa = np.dot(self.W2, hs) + self.b2 # (n_output, 1)
85     os = my_softmax(oa.T).T # (n_output, 1) softmax
86
87     cache = {"ha": ha, "hs": hs, "oa": oa, "os": os, "x": x}
88     return cache
89
90
91 def _bprop_one_example(self, cache, y):
92     """backward propagation of the network using one example.
93     Params
94         cache: a dict variable from the output of forward propagation.
95         y: int, true class index of the example
96     Returns
97         grads: gradient of the parameters of the network.
98     """
99     ha = cache["ha"] # (n_hidden, 1)
100    hs = cache["hs"] # (n_hidden, 1)
101    oa = cache["oa"] # (n_output, 1)
102    os = cache["os"] # (n_output, 1)
103    x = cache["x"] # (n_input, 1)
104    grad_oa = np.copy(os) # (n_output, 1)
105    np.put(grad_oa, y, grad_oa[y,0]-1)
106
107    grad_W2 = np.dot(grad_oa, hs.T) # (n_output, n_hidden)
108    grad_b2 = grad_oa # (n_output, 1)
109
110    grad_hs = np.dot(self.W2.T, grad_oa) # (n_hidden, 1)
111    grad_ha = np.multiply(grad_hs, np.maximum(np.sign(ha), 0))
112    # (n_hidden, 1)
113
114    grad_W1 = np.dot(grad_ha, x.T) # (n_hidden, n_input)
115    grad_b1 = grad_ha # (n_hidden, 1)
116
117    # grad_x = np.dot(self.W1.T, grad_ha) # (n_input, 1)
118    lambdas = self.lambdas
119    if not lambdas[0] == 0:
120        grad_W1 += lambdas[0] * np.sign(self.W1)

```

```

121     if not lambdas[1] == 0:
122         grad_W1 += 2 * lambdas[1] * self.W1
123     if not lambdas[2] == 0:
124         grad_W2 += lambdas[2] * np.sign(self.W2)
125     if not lambdas[3] == 0:
126         grad_W2 += 2 * lambdas[3] * self.W2
127
128     grads = {"grad_W1":grad_W1, "grad_W2":grad_W2,
129             "grad_b1":grad_b1, "grad_b2":grad_b2,
130             #"grad_hs":grad_hs, "grad_ha":grad_ha,
131             }
132     return grads
133
134
135 def fit_non_matrix_expression(self, X, y):
136     """fit X to y without using matrix expression. parameters
137     are updated after every batch_size. In each batch_size,
138     grads are accumulated by a looping over each example.
139     Params
140         X: training X ndarray (batch_size, n_features)
141         y: training labels ndarray (batch_size, )
142     Returns
143         losses: list of loss of each epoch [float]
144         errors: list of error rate of each epoch [float]
145     """
146     self._init_weights_biases(X, y)
147
148     sample_size, _ = X.shape
149     batches = int(np.ceil(sample_size / self.batch_size))
150     losses, errors = [], [] # store loss and error each epoch
151     for epoch in range(self.epochs):
152         loss, err = 0.0, 0.0
153         for j in range(batches):
154             b_start = j * self.batch_size
155             b_end = min(sample_size, (j+1) * self.batch_size)
156             batch_X = X[b_start:b_end, :]
157             batch_y = y[b_start:b_end]
158             grads = self._init_grads()
159             for k in range(b_end - b_start):
160                 cache = self._fprop_one_example(batch_X[k])
161                 grads_tmp = self._bprop_one_example(cache, batch_y[k])
162                 for key in grads:
163                     # compute mean grads[key]
164                     grads[key] += (grads_tmp[key] - grads[key])/(k + 1)
165                 loss += self._compute_one_loss(cache["os"], batch_y[k])
166             self._update_params(grads)
167         # end of mini_batch
168
169         error = self.compute_error(self.predict(X), y)
170         errors.append(error)
171         loss /= sample_size
172         loss += self._compute_regular_loss()
173         losses.append(loss)
174     # end of one epoch
175     return losses, errors
176
177
178 def _fprop(self, X):
179     """
180     forward propagation using Matrix expression
181     params

```

```

182         X: input data (batch_size, n_input)
183     returns
184         os: output of the network (batch_size, n_output)
185     """
186     ha = np.dot(X, self.W1.T) + self.b1.T # (batch_size, n_hidden)
187     hs = np.maximum(ha, 0) # relu (batch_size, n_hidden)
188
189     oa = np.dot(hs, self.W2.T) + self.b2.T # (batch_size, n_output)
190     os = my_softmax(oa) # (batch_size, n_output)
191
192     cache = {"ha": ha, "hs": hs, "oa": oa, "os": os, "X": X}
193     return cache
194
195
196 def _bprop(self, cache, y):
197     """
198     backward propagation using Matrix expression
199     params
200         os: output of network (batch_size, n_output)
201         y: true label class (batch_size, 1)
202     returns:
203         grads: gradients( grad_W1, grad_W2, grad_b1, grad_b2)
204     """
205     ha = cache["ha"] # (batch_size, n_hidden)
206     hs = cache["hs"] # (batch_size, n_hidden)
207     oa = cache["oa"] # (batch_size, n_output)
208     os = cache["os"] # (batch_size, n_output)
209     X = cache["X"] # (batch_size, n_input)
210
211     batch_size, n_o = os.shape # (batch_size, n_output)
212     grad_oa = os - one_hot(n_o, y) # (batch_size, n_output)
213     # grad_oa = np.mean(grad_oa, axis = 0, keepdims = True).T
214     # (n_output, 1)
215
216     grad_W2 = np.dot(grad_oa.T, hs) # (n_output, n_hidden)
217     grad_b2 = grad_oa # (batch_size, n_output)
218
219     grad_hs = np.dot(grad_oa, self.W2) # (batch_size, n_hidden)
220     grad_ha = np.multiply(grad_hs, np.maximum(np.sign(ha), 0))
221     # (batch_size, n_hidden)
222
223     grad_W1 = np.dot(grad_ha.T, X) # (n_hidden, n_input)
224     grad_b1 = grad_ha # (batch_size, n_hidden)
225
226     # grad_x = np.dot(grad_ha, self.W1) # (batch_size, d)
227
228     # mean grad
229     grad_W2 /= batch_size
230     grad_W1 /= batch_size
231     grad_b2 = np.mean(grad_b2, axis = 0, keepdims = True).T
232     grad_b1 = np.mean(grad_b1, axis = 0, keepdims = True).T
233
234     lambdas = self.lambdas
235     if not lambdas[0] == 0:
236         grad_W1 += lambdas[0] * np.sign(self.W1)
237     if not lambdas[1] == 0:
238         grad_W1 += 2 * lambdas[1] * self.W1
239     if not lambdas[2] == 0:
240         grad_W2 += lambdas[2] * np.sign(self.W2)
241     if not lambdas[3] == 0:
242         grad_W2 += 2 * lambdas[3] * self.W2

```



```

243
244     grads = {"grad_W1":grad_W1, "grad_W2":grad_W2,
245             "grad_b1":grad_b1, "grad_b2":grad_b2,
246             #"grad_hs":grad_hs, "grad_ha":grad_ha,
247             }
248     return grads
249
250
251 def fit(self, X, y):
252     """fit X to y. Matrix operation
253     Params
254         X: training X ndarray (sample_size, n_features)
255         y: training labels ndarray (sample_size, )
256     Returns
257         losses: list of loss of each epoch [float]
258         errors: list of error rate of each epoch [float]
259     """
260     self._init_weights_biases(X, y)
261     sample_size, _ = X.shape
262     batches = int(np.ceil(sample_size / self.batch_size))
263     losses, errors = [], []
264
265     for epoch in range(self.epochs):
266         for j in range(batches):
267             b_start = j * self.batch_size
268             b_end = min(sample_size, (j+1) * self.batch_size)
269             batch_X = X[b_start:b_end, :]
270             batch_y = y[b_start:b_end]
271             cache = self._fprop(batch_X)
272             grads = self._bprop(cache, batch_y)
273             self._update_params(grads)
274             # end batch_size
275             os = self._fprop(X)["os"]
276             loss = self.compute_empirical_risk(os, y)
277             losses.append(loss)
278
279             oy = np.argmax(os, axis = 1)
280             error = self.compute_error(oy, y)
281             errors.append(error)
282         # end epoch
283     return losses, errors
284
285
286 def predict_probs(self, X):
287     """predict probabilities (os) of given data
288     Params
289         X: examples. ndarray (certain_size, n_input)
290     Returns
291         os: probabilities of X belong to every class.
292         ndarray (certain_size, n_output)
293     """
294     return self._fprop(X)["os"]
295
296
297 def predict(self, X):
298     """predict class labels of given data
299     Params
300         X: examples ndarray (certain_size, n_input)
301     Returns
302         y_output: predicted label of given data. ndarray
303         (certain_size, )

```

```

304         """
305         os = self.predict_probs(X)
306         return np.argmax(os, axis = 1)
307
308
309     def _backup_params(self):
310         """save current network parameters
311         Params: None
312         Returns:
313             parameters saved. dict.
314         """
315         old_params = {"W1":np.copy(self.W1),
316                       "W2":np.copy(self.W2),
317                       "b1":np.copy(self.b1),
318                       "b2":np.copy(self.b2)
319                       }
320         return old_params
321
322
323     def _restore_params(self, old_params = None):
324         """restore network parameters by change network's parameters
325         Params
326             old_params: parameters restored
327         Returns:
328             None
329         """
330         if old_params is None:
331             return
332         self.W1 = old_params["W1"]
333         self.W2 = old_params["W2"]
334         self.b1 = old_params["b1"]
335         self.b2 = old_params["b2"]
336
337
338     def _update_params(self, grads, learnig_rate = None):
339         """update network parameters using given gradients and larning
340         rate
341         Params
342             grads: gradients of parameters. dict
343             learning_rate: learning rate. float
344         Returns
345             None
346         """
347         lr = learnig_rate
348         if lr is None:
349             lr = self.lr
350         self.W1 -= lr * grads["grad_W1"]
351         self.W2 -= lr * grads["grad_W2"]
352         self.b1 -= lr * grads["grad_b1"]
353         self.b2 -= lr * grads["grad_b2"]
354
355
356     def compute_error(self, oy, y):
357         """compute mis classification rate
358         Params
359             oy: output of a classifier ndarray (size, )
360             y: true class labels ndarray (size, )
361         Returns
362             error: misclassified error. float
363         """
364         return float(1 - np.sum(oy == y) / len(oy))

```

```

365
366
367 def _compute_one_loss(self, os, y):
368     """loss of one example
369     Params
370         os: column vector output (n_output, 1)
371         y: scalar int
372     Returns loss of one example. float
373     """
374     os = os.reshape(-1, 1)
375     os_y = os[y, 0] + 1e-100 # avoid np.log(0)
376     return float(-np.log(os_y))
377
378
379 def _compute_regular_loss(self):
380     """regular loss
381     Params
382         None
383     Returns
384         None
385     """
386     loss = 0.0
387     lambdas = self.lambdas
388     if not lambdas[0] == 0:
389         loss += lambdas[0] * np.sum(np.abs(self.W1))
390     if not lambdas[1] == 0:
391         loss += lambdas[1] * np.sum(np.power(self.W1, 2))
392     if not lambdas[2] == 0:
393         loss += lambdas[2] * np.sum(np.abs(self.W2))
394     if not lambdas[3] == 0:
395         loss += lambdas[3] * np.sum(np.power(self.W2, 2))
396     return float(loss)
397
398
399 def compute_empirical_risk(self, os, y):
400     """
401     compute  $\tilde{R} = \hat{R} + \mathcal{L}(\theta)$ 
402     params
403         os: output of network (sample_size, n_output)
404         y: true label (sample_size, class_index)
405     """
406     batch_size, n_output = os.shape
407     y = y.reshape(-1, )
408     y = y.astype(int)
409     os_y = os[np.arange(batch_size), y]
410     loss = np.mean(-np.log(os_y)) # might cause warning from np.log()
411     loss += self._compute_regular_loss()
412     return float(loss)

```

3. Loading two datasets: Circle dataset and Fashion MNIST dataset

Two datasets will be loaded by following codes, each of which will be shuffled and then divided into 3 parts for training, validating, and testing with a proportion about 7:2:2 and 5:1:1 respectively.

See the following codes for details.

In [3]:

```

1  from time import time
2  from datetime import datetime
3  from tqdm import tqdm
4  import utils.mnist_reader as mnist_reader
5
6  # load circles data
7  data = np.loadtxt(open('circles.txt','r'))
8  np.random.shuffle(data)
9  X = data[:,0:2]
10 y = data[:,2].reshape(-1)
11 y = y.astype(int)
12
13 total_size = len(X) # 1100
14 train_size = 700
15 valid_size = 200
16 test_size = total_size - train_size - valid_size
17
18 # split dataset to three parts: train, valid, test
19 train_X = X[0: train_size,:]
20 train_y = y[0: train_size]
21 valid_X = X[train_size: train_size + test_size,:]
22 valid_y = y[train_size: train_size + test_size]
23 test_X = X[train_size + test_size:,:]
24 test_y = y[train_size + test_size:]
25
26 debug(train_X.shape, train_y.shape)
27 debug(valid_X.shape, valid_y.shape)
28 debug(test_X.shape, test_y.shape)
29
30 # load Fashion MNIST dataset
31 fshn_train_X0, fshn_train_y0 = mnist_reader.load_mnist('data/fashion',
32                                                         kind='train')
33 fshn_test_X, fshn_test_y = mnist_reader.load_mnist('data/fashion',
34                                                      kind='t10k')
35 #plt.scatter(X[:,0], X[:,1], c = y)
36 fshn_train_X = fshn_train_X0[0:50000,:]
37 fshn_train_y = fshn_train_y0[0:50000]
38 fshn_valid_X = fshn_train_X0[50000:,:]
39 fshn_valid_y = fshn_train_y0[50000:]
40
41 debug(fshn_train_X.shape, fshn_train_y.shape)
42 debug(fshn_valid_X.shape, fshn_valid_y.shape)
43 debug(fshn_test_X.shape, fshn_test_y.shape)

```

```

(700, 2) (700,)
(200, 2) (200,)
(200, 2) (200,)
(50000, 784) (50000,)
(10000, 784) (10000,)
(10000, 784) (10000,)

```

Experiments and Results

1.

As a beginning, start with an implementation that computes the gradients for a single example, and check that the gradient is correct using the finite difference method described above.

Answer / Code

The gradient computed by forward and backward propagation relies on the methods in `MLP` class:

```
cache = mlp._fprop_one_example(x)
grads1 = mlp._bprop_one_example(cache, c)
```

The first function `_fprop_one_example(x)` receives an example `x`, reshapes it to a column vector, computes the hidden layer activation (`ha`), hidden layer output (`ho`), output layer activation (`oa`) and output layer output (`os`) in order. All these variables are cached into a dictionary `cache` as a return variable for the use of back propagation. The codes are as follows:

```
def _fprop_one_example(self, x):
    x = x.reshape(-1, 1) # to column vector
    ha = np.dot(self.W1, x) + self.b1
    hs = np.maximum(ha, 0)

    oa = np.dot(self.W2, hs) + self.b2
    os = my_softmax(oa.T).T # to row vector for softmax

    cache = {"ha": ha, "hs": hs, "oa": oa, "os": os, "x":x}
    return cache
```

The second function `mlp._bprop_one_example(cache, c)` receives the return value of the first function (`cache`) and the true class label of the example (`c`). By the equations derived in theory part and the following codes, the gradients of the network parameters are then calculated:

```

def _bprop_one_example(self, cache, y):
    ha = cache["ha"] # (n_hidden, 1)
    hs = cache["hs"] # (n_hidden, 1)
    oa = cache["oa"] # (n_output, 1)
    os = cache["os"] # (n_output, 1)
    x = cache["x"] # (n_input, 1)
    grad_oa = np.copy(os)
    np.put(grad_oa, y, grad_oa[y,0]-1)

    grad_W2 = np.dot(grad_oa, hs.T) # (n_output, n_hidden)
    grad_b2 = grad_oa # (n_output, 1)

    grad_hs = np.dot(self.W2.T, grad_oa) # (n_hidden, 1)
    grad_ha = np.multiply(grad_hs, np.maximum(np.sign(ha),0)) # (n_hidden,
1)

    grad_W1 = np.dot(grad_ha, x.T) # (n_hidden, n_input)
    grad_b1 = grad_ha # (n_hidden, 1)

    # grad_x = np.dot(self.W1.T, grad_ha) # (n_input, 1)
    lambdas = self.lambdas
    if not lambdas[0] == 0:
        grad_W1 += lambdas[0] * np.sign(self.W1)
    if not lambdas[1] == 0:
        grad_W1 += 2 * lambdas[1] * self.W1
    if not lambdas[2] == 0:
        grad_W2 += lambdas[2] * np.sign(self.W2)
    if not lambdas[3] == 0:
        grad_W2 += 2 * lambdas[3] * self.W2

    grads = {"grad_W1":grad_W1, "grad_W2":grad_W2,
            "grad_b1":grad_b1, "grad_b2":grad_b2,
            #"grad_hs":grad_hs, "grad_ha":grad_ha,
            }
    return grads

```

Two other functions `is_grad_equal` and `is_grads_equal` are implemented to check if the grads by finite difference method and by back propagation of the network are equal. Only when two values of the gradient with respect to each parameter (`W1`, `W2`, `b1`, `b2`) from two methods are almost equal (difference between 0.01), will the functions give a `True` result and print "All grads are equal".

See the following codes for details.

In [4]:

```

1  def grad_finite_diff_one_example(mlp, x, c, epsilon = 1e-5):
2      grads = mlp._init_grads()
3      params = {"grad_W1":mlp.W1, "grad_W2":mlp.W2,
4                "grad_b1":mlp.b1, "grad_b2":mlp.b2,
5                }
6      for grad_key in grads:
7          grad = grads[grad_key]
8          param = params[grad_key]
9          m, n = grad.shape
10         for i in range(m):
11             for j in range(n):
12                 old_value = param[i, j]
13
14                 value_plus = old_value + epsilon
15                 param[i, j] = value_plus
16                 cache = mlp._fprop_one_example(x)
17                 loss_plus = mlp._compute_one_loss(cache["os"], c)
18                 loss_plus += mlp._compute_regular_loss()
19
20                 value_minor = old_value - epsilon
21                 param[i, j] = value_minor
22                 cache = mlp._fprop_one_example(x)
23                 loss_minor = mlp._compute_one_loss(cache["os"], c)
24                 loss_minor += mlp._compute_regular_loss()
25
26                 grad[i, j] = (loss_plus - loss_minor)/(2.0 * epsilon)
27                 param[i, j] = old_value
28     return grads
29
30
31 def is_grad_equal(grad1, grad2, epsilon = 1e-6):
32     """check if all element values of two grad array are equal"""
33     if not grad1.shape == grad2.shape:
34         return False
35     m, n = grad1.shape
36     for i in range(m):
37         for j in range(n):
38             if abs(grad1[i,j] - grad2[i,j]) > epsilon:
39                 return False
40     return True
41
42
43 def is_grads_equal(grads1, grads2, epsilon = 1e-6):
44     for grad_key in grads1:
45         if not is_grad_equal(grads1[grad_key], grads2[grad_key], epsilon):
46             print("Failure in gradient check of param: {}".format(grad_key))
47             return False
48     print("All gradients are equal")
49     return True

```

We randomly chose an example from circle dataset, calculated gradients from two methods respectively, and checked them are equal by the following codes:

In [5]:

```
1 # create a mlp with certain neurons in hidden layer
2 mlp = MLP(n_hidden = 2)
3 mlp._init_weights_biases(X, y) # initialize weights and bias
4
5 for _ in range(1): # check times
6     index = np.random.randint(0, len(X)) # randomly select an example
7     x, c = X[index, :], y[index]
8
9     cache = mlp._fprop_one_example(x)
10    grads1 = mlp._bprop_one_example(cache, c)
11
12    grads2 = grad_finite_diff_one_example(mlp, x, c)
13    is_grads_equal(grads1, grads2)
14
```

All gradients are equal

2.

Display the gradients for both methods (direct computation and finite difference) for a small network (e.g. $d = 2$ and $d_h = 2$) with random weights and for a single example.

Answer / Code

The following codes printed the values of the gradients.

In [6]:

```

1 def print_grads(grads1, grads2, desript=["", ""]):
2     for grad_key in grads1:
3         print(grad_key + " by " + desript[0] + ":")
4         print(grads1[grad_key])
5         print(grad_key + " by " + desript[1] + ":")
6         print(grads2[grad_key])
7         print()
8
9 print_grads(grads1, grads2, ["back propagation",
10                             "finite difference method"])

```

grad_b2 by back propagation:

```

[[-0.51004331]
 [ 0.51004331]]

```

grad_b2 by finite difference method:

```

[[-0.51004331]
 [ 0.51004331]]

```

grad_W2 by back propagation:

```

[[ 0.          -0.05893403]
 [ 0.          0.05893403]]

```

grad_W2 by finite difference method:

```

[[ 0.          -0.05893403]
 [ 0.          0.05893403]]

```

grad_W1 by back propagation:

```

[[ 0.          0.          ]
 [ 0.16704075 -0.05960002]]

```

grad_W1 by finite difference method:

```

[[ 0.          0.          ]
 [ 0.16704075 -0.05960002]]

```

grad_b1 by back propagation:

```

[[0.          ]
 [0.17735494]]

```

grad_b1 by finite difference method:

```

[[0.          ]
 [0.17735494]]

```

3.

Add a hyperparameter for the minibatch size K to allow compute the gradients on a minibatch of K examples (in a matrix), by looping over the K examples (this is a small addition to your previous code).

Answer / Code

In this experiment, two member function of class `MLP` were used:

```

_compute_one_loss(self, os, y)
_compute_regular_loss(self)

```

The loss of one example is calculated by the following equation:

$$L(\mathbf{x}, y) = -\log \mathbf{o}_y^s(\mathbf{x})$$

which we have already seen in question 2.4 of theoretical part. The corresponding codes are:

```
def _compute_one_loss(self, os, y):
    """loss of one example
    params
        os: column vector output (n_output, 1)
        y: scalar int
    """
    os = os.reshape(-1, 1)
    os_y = os[y, 0] + 1e-100 # avoid np.log(0)
    return float(-np.log(os_y))
```

The loss of regularization is also a part of the loss on either one example or batch examples, and they can be derived by following codes:

```
def _compute_regular_loss(self):
    loss = 0.0
    lambdas = self.lambdas
    if not lambdas[0] == 0:
        loss += lambdas[0] * np.sum(np.abs(self.W1))
    if not lambdas[1] == 0:
        loss += lambdas[1] * np.sum(np.power(self.W1, 2))
    if not lambdas[2] == 0:
        loss += lambdas[2] * np.sum(np.abs(self.W2))
    if not lambdas[3] == 0:
        loss += lambdas[3] * np.sum(np.power(self.W2, 2))
    return float(loss)
```

For the method of back propagation, the gradient with respect to a parameter on a batch size of K examples can be considered as the mean of the gradient on each sample in the batch, which is implemented in the function: `grad_back_prop_batch(mlp, X, y)`.

For the method of finite difference method, the gradient relies on the tiny change of average loss by the tiny change of the parameter. So, the key is to compute the average loss on the minibatch. See the implementation of the function `grad_finite_diff_batch(mlp, X, y, epsilon = 1e-5)` for details.

In [7]:

```

1  def grad_back_prop_batch(mlp, X, y):
2      batch_size, n_input = X.shape
3      y = y.reshape(-1, )
4      grads1 = mlp._init_grads()
5      for i in range(batch_size):
6          x, c = X[i,:], y[i]
7          cache = mlp._fprop_one_example(x)
8          grads = mlp._bprop_one_example(cache, c)
9          for key in grads1:
10             # compute mean grads1[key]
11             grads1[key] += (grads[key] - grads1[key])/(i + 1)
12     return grads1
13
14
15 def grad_finite_diff_batch(mlp, X, y, epsilon = 1e-5):
16     grads = mlp._init_grads()
17     params = {"grad_W1":mlp.W1, "grad_W2":mlp.W2,
18              "grad_b1":mlp.b1, "grad_b2":mlp.b2,
19              }
20     batch_size, _ = X.shape
21     y = y.reshape(-1, )
22     y = y.astype(int)
23     for grad_key in grads:
24         grad = grads[grad_key]
25         param = params[grad_key]
26         m, n = grad.shape
27         for i in range(m):
28             for j in range(n):
29                 old_value = param[i, j]
30
31                 value_plus = old_value + epsilon
32                 param[i, j] = value_plus
33                 loss_plus = 0
34                 for k in range(batch_size):
35                     cache = mlp._fprop_one_example(X[k,:])
36                     loss_plus += mlp._compute_one_loss(cache["os"], y[k])
37                 loss_plus /= batch_size
38                 loss_plus += mlp._compute_regular_loss()
39
40                 value_minor = old_value - epsilon
41                 param[i, j] = value_minor
42                 loss_minor = 0
43                 for k in range(batch_size):
44                     cache = mlp._fprop_one_example(X[k,:])
45                     loss_minor += mlp._compute_one_loss(cache["os"], y[k])
46                 loss_minor /= batch_size
47                 loss_minor += mlp._compute_regular_loss()
48                 grad[i, j] = (loss_plus - loss_minor)/(2.0 * epsilon)
49                 param[i, j] = old_value
50
51     return grads

```

The graients from two methods should be equal. Again, we can check them by the functions used in the experiment 1.

Luckly, We got the same gradients.

In [8]:

```
1 # create a mlp with certain neurons in hidden layer
2 mlp = MLP(n_hidden = 2)
3 mlp._init_weights_biases(X, y) # initialize weights and bias
4 batch_size = 10
5 start = np.random.randint(0, len(X)-batch_size)
6 batch_X = X[start:start + batch_size,:]
7 batch_y = y[start:start + batch_size]
8
9 grads1 = grad_back_prop_batch(mlp, batch_X, batch_y)
10 grads2 = grad_finite_diff_batch(mlp, batch_X, batch_y)
11 is_grads_equal(grads1, grads2)
```

All gradients are equal

Out[8]:

True

4.

Display the gradients for both methods (direct computation and finite difference) for a small network (e.g. $d = 2$ and $d_h = 2$) with random weights and for a minibatch with 10 examples (you can use examples from both classes from the two circles dataset).

Answers / Code

Just use the function implemented in experiment 2 to display the two gradients.

In [9]:

```
1 print_grads(grads1, grads2, ["back propagation", "finite difference method"])
```

grad_b2 by back propagation:

```
[[-0.2040551]
 [ 0.2040551]]
```

grad_b2 by finite difference method:

```
[[-0.2040551]
 [ 0.2040551]]
```

grad_W2 by back propagation:

```
[[ 0.0170263 -0.00454214]
 [-0.0170263  0.00454214]]
```

grad_W2 by finite difference method:

```
[[ 0.0170263 -0.00454214]
 [-0.0170263  0.00454214]]
```

grad_W1 by back propagation:

```
[[-9.75817411e-05 -1.08023074e-03]
 [ 2.85217777e-02 -4.04200060e-02]]
```

grad_W1 by finite difference method:

```
[[-9.75817427e-05 -1.08023074e-03]
 [ 2.85217777e-02 -4.04200060e-02]]
```

grad_b1 by back propagation:

```
[[-0.0009252 ]
 [ 0.01873617]]
```

grad_b1 by finite difference method:

```
[[-0.0009252 ]
 [ 0.01873617]]
```

5.

Train your neural network using gradient descent on the two circles dataset. Plot the decision regions for several different values of the hyperparameters (weight decay, number of hidden units, early stopping) so as to illustrate their effect on the capacity of the model.

Answer / Code

We specifically implemented the following method to train our network on different hyperparameters and plot the corresponding decision regions.

As to early stopping, we set the tolerant epoch times `n_patience = 20`. When error rate on validate dataset starts to increase, we observe the consequent `n_patience` number of epochs before stopping training.

Several member functions of class `MLP` are used in this experiments:

```
_init_weights_biases(train_X, train_y)
_fprop_one_example(batch_X[k])
_bprop_one_example(cache, batch_y[k])
_update_params(grads)
_compute_regular_loss()
predict_probs(train_X)
compute_error(train_oy, train_y)
compute_empirical_risk(valid_os, valid_y)
_backup_params()
_restore_params((best_params))
```

The purposes of these functions are just as simple as the names of the functions tell, and they are also easy to implement. No further explanation for these functions. See the implementation of the class `MLP` for details.

In [10]:

```

1  def train_by_looping(train_X, train_y, valid_X, valid_y,
2      n_hidden = 8, early_stop = False, weight_decay = [0.0, 0.0, 0.0, 0.0]
3      learning_rate = 1e-3, epochs = 1000, batch_size = 10
4      ):
5      mlp = MLP(n_hidden = n_hidden,
6          learning_rate = learning_rate,
7          epochs = epochs,
8          batch_size = batch_size,
9          early_stop = early_stop,
10         lambdas = weight_decay)
11
12     mlp._init_weights_biases(train_X, train_y)
13
14     # for early stopping
15     best_params = None # params that performs best in valid set
16     n_patience = 200 # maximum numbers of epoch tolerated when mlp
17         # with current params doesn't perform better than
18         # with best_params.
19     n_epochs_bad = 0 # epochs which passed yet didn't bring better params
20     best_valid_loss = float('inf')
21     best_valid_error = float('inf')
22
23
24     sample_size, _ = train_X.shape
25     batches = int(np.ceil(sample_size / mlp.batch_size))
26     losses = [[], []] # losses[0] is for training, 1 for valid
27     errors = [[], []] # [0] for training, 1 for valid
28
29     for epoch in range(mlp.epochs):
30         train_loss = 0.0
31         for j in range(batches):
32             b_start = j * mlp.batch_size
33             b_end = min(sample_size, (j+1) * mlp.batch_size)
34             batch_X = X[b_start:b_end, :]
35             batch_y = y[b_start:b_end]
36             grads = mlp._init_grads()
37             for k in range(b_end - b_start):
38                 cache = mlp._fprop_one_example(batch_X[k])
39                 grads_tmp = mlp._bprop_one_example(cache, batch_y[k])
40                 for key in grads:
41                     # compute mean grads[key]
42                     grads[key] += (grads_tmp[key] - grads[key]) / (k + 1)
43                 train_loss += mlp._compute_one_loss(cache["os"], batch_y[k])
44             mlp._update_params(grads)
45         # end of mini_batch
46
47         train_loss /= sample_size
48         train_loss += mlp._compute_regular_loss()
49         losses[0].append(train_loss)
50
51         if early_stop == True:
52             train_os = mlp.predict_probs(train_X)
53             train_oy = np.argmax(train_os, axis = 1)
54             train_error = mlp.compute_error(train_oy, train_y)
55             errors[0].append(train_error)
56
57             valid_os = mlp.predict_probs(valid_X)
58             valid_oy = np.argmax(valid_os, axis = 1)
59             valid_loss = mlp.compute_empirical_risk(valid_os, valid_y)

```

```
60     losses[1].append(valid_loss)
61     valid_error = mlp.compute_error(valid_oy, valid_y)
62     errors[1].append(valid_error)
63
64     if epoch == 0: # first epoch completed
65         best_params = mlp._backup_params()
66         best_valid_error = valid_error
67     else:
68         if valid_error > best_valid_error: #
69             n_epochs_bad += 1
70             if n_epochs_bad == n_patientce: # stop
71                 mlp._restore_params(best_params)
72                 return mlp, losses, errors, best_valid_error
73             else: # loss is decreasing
74                 best_params = mlp._backup_params()
75                 best_valid_error = valid_error
76                 n_epochs_bad = 0
77     # end of one epoch
78     return mlp, losses, errors, best_valid_error
```

We trained the network using different combination of the following values of the hyperparameters.

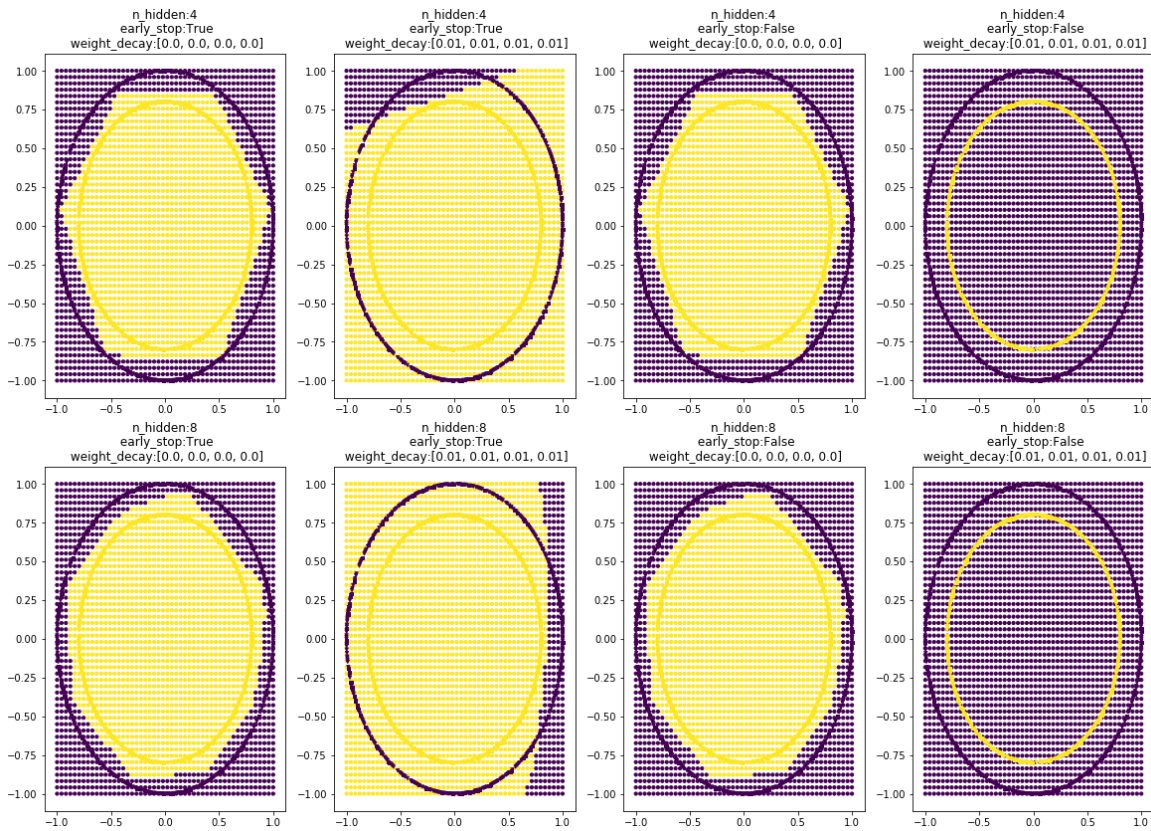
1. number of hidden layers: 4 and 8
2. early stop: True and False
3. weight decays(λ s): [0.0, 0.0, 0.0, 0.0] and [0.01, 0.01, 0.01, 0.01]

In [11]:

```

1  # hyper parameters selected.
2  n_hiddens = [4, 8]
3  early_stops = [True, False]
4  weight_decays = [[0.0]*4, [0.01]*4]
5
6  train_valid_X = np.vstack((train_X, valid_X))
7  train_valid_y = np.vstack((train_y.reshape(-1,1), valid_y.reshape(-1,1)))
8  train_valid = np.hstack((train_valid_X, train_valid_y))
9
10 min_x1, max_x1 = min(train_valid[:,0]), max(train_valid[:,0])
11 min_x2, max_x2 = min(train_valid[:,1]), max(train_valid[:,1])
12
13 n_points = 50
14 grid_x1 = np.linspace(min_x1, max_x1, num = n_points)
15 grid_x2 = np.linspace(min_x2, max_x2, num = n_points)
16
17 X1, X2 = np.meshgrid(grid_x1, grid_x2)
18 grid_data = np.hstack((X1.reshape(-1, 1), X2.reshape(-1, 1)))
19 len_i = len(n_hiddens)
20 len_j = len(early_stops)
21 len_k = len(weight_decays)
22
23 plt.figure(figsize=(20, 14))
24 for i in range(len_i):
25     for j in range(len_j):
26         for k in range(len_k):
27             mlp, _, _, _ = train_by_looping(train_X, train_y,
28                                             valid_X, valid_y,
29                                             n_hiddens[i],
30                                             early_stops[j],
31                                             weight_decays[k],
32                                             epochs = 2000
33                                             )
34             predicts = mlp.predict(grid_data)
35             plt.subplot(len_i, len_j * len_k,
36                         i * (len_j * len_k) + (j * len_k) + k + 1)
37             plt.scatter(grid_data[:,0], grid_data[:,1], c = predicts, s=10)
38             #The training points
39             plt.scatter(train_X[:,0], train_X[:,1], c = train_y,
40                         marker = 'v', s=10)
41             # The test points
42             plt.scatter(valid_X[:,0], valid_X[:,1], c = valid_y,
43                         marker = 's', s=10)
44             plt.title("n_hidden:{}\n early_stop:{}\n weight_decay:{}"
45                      .format(n_hiddens[i], early_stops[j], weight_decays[k]))
46             # debug("{}".format(i * (len_j * len_k) + (j * len_k) + k + 1))
47 #plt.subplots_adjust()
48 plt.show()

```



6.

As a second step, copy your existing implementation to modify it to a new implementation that will use matrix calculus (instead of a loop) on batches of size K to improve efficiency. **Take the matrix expressions in numpy derived in the first part, and adapt them for a minibatch of size K . Show in your report what you have modified (describe the former and new expressions with the shapes of each matrices).**

Answer / Code

As we have seen before, forward and backward propagation computation are implemented using loop over single example in the following two functions:

```
_fprop_one_example(self, x)
_bprop_one_example(self, cache, y)
```

For this experiment, the above functions are accordingly implemented in matrix expressions in:

```
_fprop(self, X)
_bprop(self, cache, y)
```

The matrix expression are shown as follows with comments on each line shows the shape of the corresponding variable.

forward propagation on one example

```
x = x.reshape(-1, 1)           # (n_input, 1)
ha = np.dot(self.W1, x) + self.b1 # (n_hidden, 1)
hs = np.maximum(ha, 0)          # (n_hidden, 1)
oa = np.dot(self.W2, hs) + self.b2 # (n_output, 1)
os = my_softmax(oa.T).T         # (n_output, 1)
```

backward propagation on one example

```

grad_oa = np.copy(os) # (n_output, 1)
np.put(grad_oa, y, grad_oa[y,0]-1)
grad_W2 = np.dot(grad_oa, hs.T) # (n_output, n_hidden)
grad_b2 = grad_oa # (n_output, 1) # (n_output, 1)
grad_hs = np.dot(self.W2.T, grad_oa) # (n_hidden, 1)
grad_ha = np.multiply(grad_hs, np.maximum(np.sign(ha), 0)) # (n_hidden, 1)
grad_W1 = np.dot(grad_ha, x.T) # (n_hidden, n_input)
grad_b1 = grad_ha # (n_hidden, 1)

```

forward propagation on batch examples

```

ha = np.dot(X, self.W1.T) + self.b1.T # (batch_size, n_hidden)
hs = np.maximum(ha, 0) # relu # (batch_size, n_hidden)
oa = np.dot(hs, self.W2.T) + self.b2.T # (batch_size, n_output)
os = my_softmax(oa) # (batch_size, n_output)

```

backward propagation on batch examples

```

batch_size, n_o = os.shape # (batch_size, n_output)
grad_oa = os - one_hot(n_o, y) # (batch_size, n_output)
# grad_oa = np.mean(grad_oa, axis = 0, keepdims = True).T
# (n_output, 1)

grad_W2 = np.dot(grad_oa.T, hs) # (n_output, n_hidden)
grad_b2 = grad_oa # (batch_size, n_output)
grad_hs = np.dot(grad_oa, self.W2) # (batch_size, n_hidden)
grad_ha = np.multiply(grad_hs, np.maximum(np.sign(ha), 0))
# (batch_size, n_hidden)

grad_W1 = np.dot(grad_ha.T, X) # (n_hidden, n_input)
grad_b1 = grad_ha # (batch_size, n_hidden)
# grad_x = np.dot(grad_ha, self.W1) # (batch_size, n_input)

# mean grad
grad_W2 /= batch_size # shape keeps
grad_W1 /= batch_size # shape keeps
grad_b2 = np.mean(grad_b2, axis = 0, keepdims = True).T # (n_output, 1)
grad_b1 = np.mean(grad_b1, axis = 0, keepdims = True).T # (n_hidden, 1)

```

For both approaches, the shapes of the parameters (W1, W2, b1, b2) do not change; they remain the same shape when they are initialized:

```

# W1 (n_hidden, n_input)
# W2 (n_output, n_hidden)
# b1 (n_hidden, 1)
# b2 (n_output, 1)

```

See the implementation of class `MLP` for the complete implementation of **forward propagation** and **backward propagation** on both one example and batch_size examples

7.

Compare both implementations (with a loop and with matrix calculus) to check that they both give the same values for the gradients on the parameters, first for $K = 1$, then for $K = 10$. Display the gradients for both methods.

Answer / Code

In this experiment, we use the function `grad_back_prop_batch(mlp, batch_X, batch_y)` implemented in **experiment 3** to calculate the gradient using loop on `batch_size` (K) examples for both $K = 1$ and $K = 10$, and use the functions:

```
_fprop(batch_X)
_bprop(cache, batch_y)
```

to calculate the gradient using matrix expression on `batch_size` (K) examples for both $K = 1$ and $K = 10$.

Here is the code and results:

In [12]:

```

1 Ks = [1, 10]
2 for K in Ks:
3     print("Gradient check by loop and matrix on batch_size(K) = {}".format(K))
4     mlp = MLP(n_hidden = 2, batch_size = K)
5     mlp._init_weights_biases(X, y) # initialize weights and bias
6     batch_size = K
7     start = np.random.randint(0, len(X)-batch_size)
8     batch_X = X[start:start + batch_size,:]
9     batch_y = y[start:start + batch_size]
10
11     grads_loop = grad_back_prop_batch(mlp, batch_X, batch_y)
12
13     cache = mlp._fprop(batch_X)
14     grads_matrix = mlp._bprop(cache, batch_y)
15
16     is_grads_equal(grads_loop, grads_matrix)
17
18     print_grads(grads_loop, grads_matrix, ["using loop", "matrix operation"])

```

Gradient check by loop and matrix on batch_size(K) = 1

All gradients are equal

grad_b2 by using loop:

```

[[-0.51004331]
 [ 0.51004331]]

```

grad_b2 by matrix operation:

```

[[-0.51004331]
 [ 0.51004331]]

```

grad_W2 by using loop:

```

[[ 0.          -0.05893403]
 [ 0.           0.05893403]]

```

grad_W2 by matrix operation:

```

[[ 0.          -0.05893403]
 [ 0.           0.05893403]]

```

grad_W1 by using loop:

```

[[ 0.          0.          ]
 [ 0.16704075 -0.05960002]]

```

grad_W1 by matrix operation:

```

[[ 0.          0.          ]
 [ 0.16704075 -0.05960002]]

```

grad_b1 by using loop:

```

[[0.          ]
 [0.17735494]]

```

grad_b1 by matrix operation:

```

[[0.          ]
 [0.17735494]]

```

Gradient check by loop and matrix on batch_size(K) = 10

All gradients are equal

grad_b2 by using loop:

```

[[-0.2040551]
 [ 0.2040551]]

```

grad_b2 by matrix operation:

```

[[-0.2040551]
 [ 0.2040551]]

```

grad_W2 by using loop:

```
[[ 0.0170263 -0.00454214]
 [-0.0170263  0.00454214]]
grad_W2 by matrix operation:
[[ 0.0170263 -0.00454214]
 [-0.0170263  0.00454214]]

grad_W1 by using loop:
[[-9.75817411e-05 -1.08023074e-03]
 [ 2.85217777e-02 -4.04200060e-02]]
grad_W1 by matrix operation:
[[-9.75817411e-05 -1.08023074e-03]
 [ 2.85217777e-02 -4.04200060e-02]]

grad_b1 by using loop:
[[-0.0009252 ]
 [ 0.01873617]]
grad_b1 by matrix operation:
[[-0.0009252 ]
 [ 0.01873617]]
```

8.

Time how long takes an epoch on fashion MNIST (1 epoch = 1 full traversal through the whole training set) for $K = 100$ for both versions (loop over a minibatch and matrix calculus).

Answer / Code

In this experiment, we use two member functions in our class `MLP` :

```
def fit_non_matrix_expression(self, X, y):
def fit(self, X, y)
```

Both functions using the hyper parameters direvied by function `__init__()` . Main parameters are:

```
n_hidden = 2,
lambdas = [0, 0, 0, 0],
learning_rate = 1e-3,
epochs = 200,
batch_size = 10,
early_stop = False
```

The first function tries to fit the dataset (X, y) by looping over each example, whereas the second does same thing by matrix calculus. See implementation of `MLP` for details of two functions.

To see how long **one** epoch using **batch_size** ($K = 100$) would take on two training methods, we first created an instance `m1p` by the code:

```
K = 100
m1p = MLP(n_hidden = 100, batch_size = K, epochs = 1)
```

Then, by execute the two function, we could calculate the time elapsed on each method.

Here is the complete codes and results:

In [13]:

```

1 K = 100
2 mlp = MLP(n_hidden = 100, batch_size = K, epochs = 1)
3 time_start = time()
4 mlp.fit_non_matrix_expression(fshn_train_X, fshn_train_y)
5 print("{:.2f} seconds for 1 epoch using loop over mini_batch (K={})."
6       .format(time() - time_start, K))
7
8 time_start = time()
9 mlp.fit(fshn_train_X, fshn_train_y)
10 print("{:.2f} seconds for 1 epoch using matrix calculus on mini_batch (K={})."
11       .format(time() - time_start, K))

```

18.67 seconds for 1 epoch using loop over mini_batch (K=100).

1.08 seconds for 1 epoch using matrix calculus on mini_batch (K=100).

9.

Adapt your code to compute the error (proportion of misclassified examples) on the training set as well as the total loss on the training set during each epoch of the training procedure, and at the end of each epoch, it computes the error and average loss on the validation set and the test set. Display the 6 corresponding figures (error and average loss on train/valid/test), and write them in a log file.

Answer / Code

Two functions are implemented perform this experiment:

1. `train` : train the network on training set, validation set, and test set; compute loss and error on three sets in each epoch; write them into a log file. Early stopping may also used in this experiment. In this function, loss and error on each epoch are calculated with the help of three member functions of class:

```

predict_probs(X)
compute_empirical_risk(os, y)
compute_error(oy, y)

```

See the class implementation for details.

2. `plot_learning_curves` : plot loss and error curves on the three data set.

In [14]:

```

1  def train(train_X, train_y, valid_X, valid_y, test_X, test_y,
2          n_hidden = 100, early_stop = False, weight_decay = [0.0]*4,
3          learning_rate = 1e-3, epochs = 2000,
4          batch_size = 64, n_patience = 10,
5          log_file_name = "mlp_training_log"):
6
7      mlp = MLP(n_hidden = n_hidden,
8               learning_rate = learning_rate,
9               epochs = epochs,
10              batch_size = batch_size,
11              early_stop = early_stop,
12              lambdas = weight_decay)
13
14      mlp._init_weights_biases(train_X, train_y)
15      # for early stopping
16      best_params = None # params that performs best in valid set
17      # n_patience = 20 # maximum numbers of epoch tolerated when mlp
18      # with current params doesn't perform better than
19      # with best_params.
20      n_epochs_bad = 0 # epochs which passed yet didn't bring better params
21      best_valid_error = float('inf')
22
23      sample_size, _ = train_X.shape
24      batches = int(np.ceil(sample_size / mlp.batch_size))
25      losses = [[], [], []] # [[train], [valid], [test]]
26      errors = [[], [], []] # [[train], [valid], [test]]
27
28      ISOTIMEFORMAT = '%Y-%m-%d %H:%M:%S'
29      file_name = log_file_name + str(datetime.now().strftime(ISOTIMEFORMAT))
30      file_name += ".txt"
31      f = open(file_name, "w")
32      f.write("training records of a MLP network\n")
33      f.close()
34      for epoch in range(mlp.epochs):
35          # for epoch in range(mlp.epochs):
36              # batch matrix forward and backward
37              for j in range(batches):
38                  batch_start = j * mlp.batch_size
39                  batch_end = min(sample_size, (j+1) * mlp.batch_size)
40                  batch_X = train_X[batch_start:batch_end, :]
41                  batch_y = train_y[batch_start:batch_end]
42                  cache = mlp._fprop(batch_X)
43                  grads = mlp._bprop(cache, batch_y)
44                  mlp._update_params(grads)
45              # end of mini_batch
46
47              # compute loss and error rate on train, valid, and test dataset
48              dataset = [(train_X, train_y), (valid_X, valid_y), (test_X, test_y)]
49              for i in range(len(dataset)):
50                  os = mlp.predict_probs(dataset[i][0])
51                  loss = mlp.compute_empirical_risk(os, dataset[i][1])
52                  losses[i].append(loss)
53                  oy = np.argmax(os, axis = 1)
54                  error = mlp.compute_error(oy, dataset[i][1])
55                  errors[i].append(error)
56
57              # write to log file for current epoch
58              f = open(file_name, "a+")
59              log = "epoch {:<5}: ".format(epoch)

```



```

60     log += "train_l:{:5.3f}, train_er:{:.3f}; ".format(
61         losses[0][-1], errors[0][-1])
62     log += "valid_l:{:5.3f}, valid_er:{:.3f}; ".format(
63         losses[1][-1], errors[1][-1])
64     log += "test_l:{:5.3f}, test_er:{:.3f}; ".format(
65         losses[2][-1], errors[2][-1])
66     log += "\n"
67     f.write(log)
68     f.close()
69
70     # early stopping check and related process
71     if early_stop == True:
72         if epoch == 0: # first epoch completed
73             best_params = mlp._backup_params()
74             best_valid_error = errors[1][-1] # using error
75         else:
76             latest_valid_error = errors[1][-1] # using error
77             if latest_valid_error > best_valid_error: # use > instead of >
78                 n_epochs_bad += 1
79                 if n_epochs_bad == n_patience: # stop
80                     mlp._restore_params(best_params)
81                     best_epoch = epoch - n_patience
82                     f = open(file_name, "a+")
83                     log = "Early stop at epoch {}. With".format(best_epoch)
84                     log += " valid_loss:{:.3f}, valid_error:{:.3f}".format(
85                         losses[1][best_epoch], errors[1][best_epoch])
86                     log += " train_loss:{:.3f}, train_error:{:.3f}".format(
87                         losses[0][best_epoch], errors[0][best_epoch])
88                     f.write(log)
89                     f.close()
90                     debug(log)
91                     debug("log file: {}".format(file_name))
92                     return mlp, losses, errors, best_valid_error
93             else: # loss is decreasing
94                 best_params = mlp._backup_params()
95                 best_valid_error = latest_valid_error
96                 n_epochs_bad = 0
97     # end of one epoch
98     debug("training complete. log file: {}".format(file_name))
99     f = open(file_name, "a+")
100     log = "n_hidden:{:<5}\t".format(n_hidden)
101     log += "learning_rate:{:.3f}\t".format(learning_rate)
102     log += "batch_size:{:<3}\t".format(batch_size)
103     log += "n_patience:{:<3}\t".format(n_patience)
104     log += "weight_decay:{}\n".format(weight_decay)
105     f.write(log)
106     f.close()
107     return mlp, losses, errors, best_valid_error
108
109
110 def plot_learning_curves(losses, errors, data_set_name = ""):
111     support_x = np.arange(len(losses[0]))
112     colors = ['g-', 'b-', 'r']
113     set_groups = ["train_set", "valid_set", "test_set"]
114
115     plt.figure(figsize=(13, 6))
116     plt.grid(True) # add a grid
117
118     plt.subplot(1, 2, 1)
119     plt.title("Loss curves on {} dataset".format(data_set_name))
120     legends = []

```

```

121     for i in range(3):
122         plt.xlabel("epoch")
123         plt.ylabel("loss")
124         plt.plot(support_x, losses[i], colors[i])
125         legends.append(set_groups[i])
126     plt.legend(legends)
127
128     plt.subplot(1, 2, 2)
129     plt.title("Error curves on {} dataset".format(data_set_name))
130     for i in range(3):
131         plt.xlabel("epoch")
132         plt.ylabel("error")
133         plt.plot(support_x, errors[i], colors[i])
134     plt.legend(legends)
135     plt.show()

```

The following codes create an MLP instance with the hyper parameters given, train the network, give the best error rate on validation set, and plot the loss and error curves.

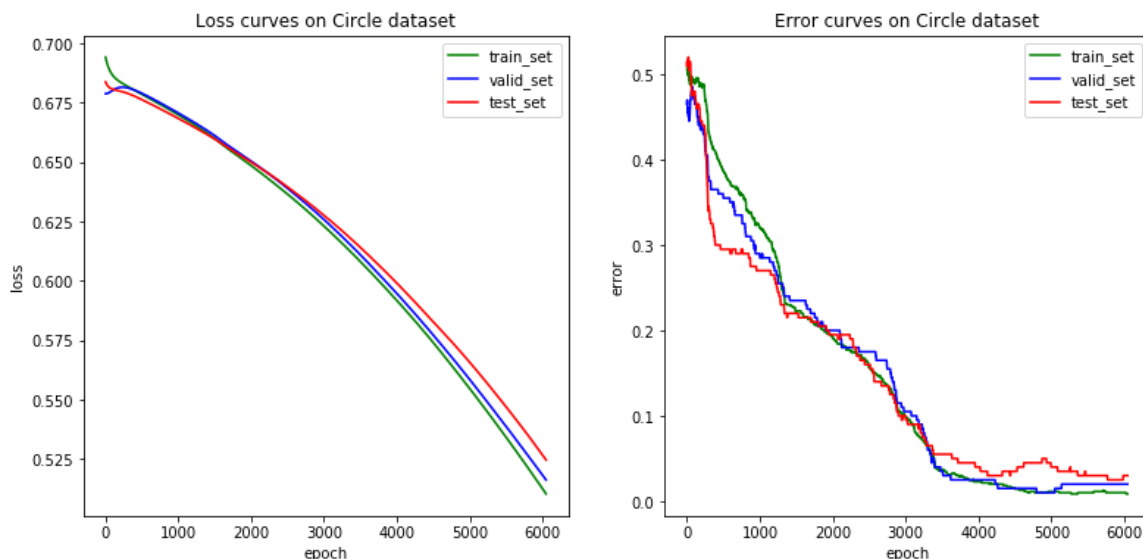
In [15]:

```

1  n_hidden, early_stop, weight_decay = 16, True, [0.0]*4
2  mlp, losses, errors, best_valid_error = train(train_X, train_y,
3                                              valid_X, valid_y,
4                                              test_X, test_y,
5                                              n_hidden = n_hidden,
6                                              early_stop = early_stop,
7                                              weight_decay = weight_decay,
8                                              learning_rate = 1e-3,
9                                              epochs = 8000,
10                                             batch_size = 64,
11                                             n_patience = 1000,
12                                             log_file_name = "mlp_circle")
13
14  debug("best_valid_error: ", best_valid_error)
15  plot_learning_curves(losses, errors, "Circle")

```

Early stop at epoch 5048. With valid_loss:0.557, valid_error:0.010 train_loss:0.553, train_error:0.011
 log file: mlp_circle2018-11-09 08:45:55.txt
 best_valid_error: 0.010000000000000009



10

Train your network on the fashion MNIST dataset. Plot the training/valid/test curves (error and loss as a function of the epoch number, corresponding to what you wrote in a file in the last question). Add to your report the curves obtained using your best hyperparameters, i.e. for which you obtained your best error on the validation set. We suggest 2 plots : the first one will plot the error rate (train/valid/test with different colors, show which color in a legend) and the other one for the averaged loss (on train/valid/test). You should be able to get less than 20% test error.

Answer / Code

As for fashion MNIST dataset, function `train` implemented in last experiment is used. This time, the hyper parameters are set as follows:

```
n_hidden = 128
early_stop = True
weight_decay = [0.0, 0.0, 0.0, 0.0]
learning_rate = 1e-4
batch_size = 128
```

The best error rate on validation set is: 12.9%.

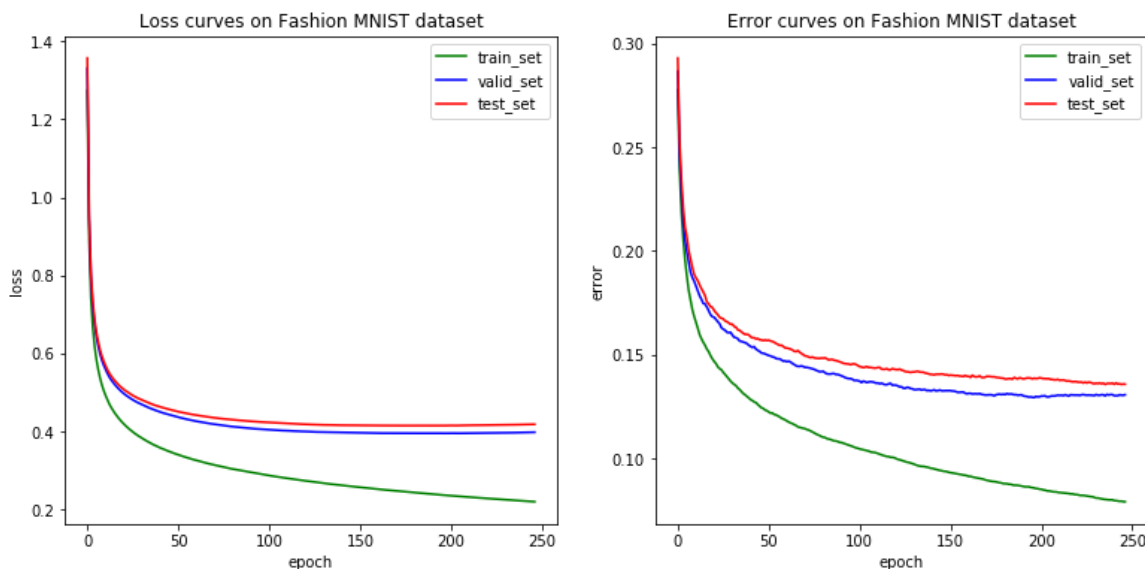
In [16]:

```

1 mlp, losses, errors, best_valid_error = train(fshn_train_X, fshn_train_y,
2       fshn_valid_X, fshn_valid_y,
3       fshn_test_X, fshn_test_y,
4       n_hidden = 128,
5       early_stop = True,
6       weight_decay = [0.0]*4,
7       learning_rate = 1e-4,
8       epochs = 1000,
9       batch_size = 128,
10      n_patience = 50,
11      log_file_name = "mlp_fashion_mnist")
12
13 debug("best_valid_error: {:.3f}".format(best_valid_error))
14 plot_learning_curves(losses, errors, "Fashion MNIST")

```

Early stop at epoch 196. With valid_loss:0.395, valid_error:0.129 tra
in_loss:0.237, train_error:0.086
log file: mlp_fashion_mnist2018-11-09 08:46:07.txt
best_valid_error: 0.129



The following codes are used to select the best from our hyper-parameter set:

```

1 n_hiddens = [64, 128]
2 early_stops = [True]
3 weight_decays = [[0.0]*4, [0, 0.001, 0.0, 0.001], [0.01]*4, [0.1]*4]
4 learning_rates = [1e-3, 1e-4]
5 batch_sizes = [64, 128]
6
7 for n_hidden in n_hiddens:
8     for early_stop in early_stops:
9         for weight_decay in weight_decays:
10             for learning_rate in learning_rates:
11                 for batch_size in batch_sizes:
12                     mlp, losses, errors, best_valid_error = train(
13                         fshn_train_X, fshn_train_y,
14                         fshn_valid_X, fshn_valid_y,
15                         fshn_test_X, fshn_test_y,
16                         n_hidden = n_hidden,
17                         early_stop = early_stop,
18                         weight_decay = weight_decay,

```

```
19         learning_rate = learning_rate,
20         epochs = 1000,
21         batch_size = batch_size,
22         n_patience = 20,
23         log_file_name = "mlp_fashion_mnist")
24     log = "best_v_error:{:.3f}\t".format(best_valid_error)
25     log += "n_hidden:{:<4}\t".format(n_hidden)
26     log += "early_stop:{}\t".format(early_stop)
27     log += "learning_rate:{:.4f}\t".format(learning_rate)
28     log += "batch_size:{:<4}\t".format(batch_size)
29     log += "lambdas:{}\n".format(weight_decay)
30     debug(log)
```

The end of the report of homework3 IFT6390

Team Member

Qiang Ye (20139927), Lifeng Wan (20108546)

Coding Environment

python 3.5.2, numpy 1.14.2 matplotlib 2.2.0

In []:

1