

## Lab 2: Classification

### Machine Learning 1, September 2016

- The lab exercises should be made in groups of two people.
- The deadline is October 9th (Sunday) 23:59.
- Assignment should be sent to your teaching assistant. The subject line of your email should be "lab#\_lastname1\_lastname2\_lastname3".
- Put your and your teammates' names in the body of the email.
- Attach the .IPYNB (IPython Notebook) file containing your code and answers. Naming of the file follows the same rule as the subject line. For example, if the subject line is "lab01\_Kingma\_Hu", the attached file should be "lab01\_Kingma\_Hu.ipynb". Only use underscores ("\_") to connect names, otherwise the files cannot be parsed.

Notes on implementation:

- For this notebook you need to answer a few theory questions, add them in the Markdown cell's below the question. Note: you can use Latex-style code in here.
- Focus on Part 1 the first week, and Part 2 the second week!
- You should write your code and answers below the questions in this IPython Notebook.
- Among the first lines of your notebook should be "%pylab inline". This imports all required modules, and your plots will appear inline.
- If you have questions outside of the labs, post them on blackboard or email me.
- NOTE: Make sure we can run your notebook / scripts!

```
In [1]: %pylab inline
import gzip, cPickle
```

Populating the interactive namespace from numpy and matplotlib

## Part 1. Multiclass logistic regression

Scenario: you have a friend with one big problem: she's completely blind. You decided to help her: she has a special smartphone for blind people, and you are going to develop a mobile phone app that can do *machine vision* using the mobile camera: converting a picture (from the camera) to the meaning of the image. You decide to start with an app that can read handwritten digits, i.e. convert an image of handwritten digits to text (e.g. it would enable her to read precious handwritten phone numbers).

A key building block for such an app would be a function `predict_digit(x)` that returns the digit class of an image patch `x`. Since hand-coding this function is highly non-trivial, you decide to solve this problem using machine learning, such that the internal parameters of this function are automatically learned using machine learning techniques.

The dataset you're going to use for this is the MNIST handwritten digits dataset (<http://yann.lecun.com/exdb/mnist/>). You can load the data from `mnist.pkl.gz` we provided, using:

```
In [2]: def load_mnist():
        f = gzip.open('mnist.pkl.gz', 'rb')
        data = cPickle.load(f)
        f.close()
        return data

(x_train, t_train), (x_valid, t_valid), (x_test, t_test) = load_mnist()
```

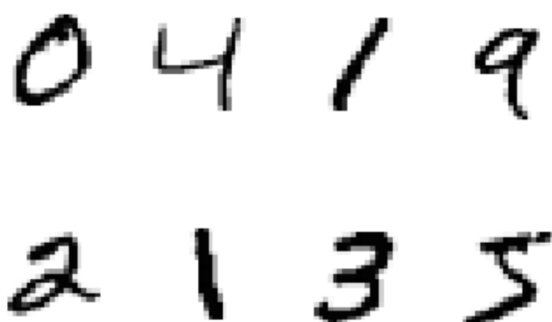
The tuples represent train, validation and test sets. The first element ( $x_{\text{train}}$ ,  $x_{\text{valid}}$ ,  $x_{\text{test}}$ ) of each tuple is a  $N \times M$  matrix, where  $N$  is the number of datapoints and  $M = 28^2 = 784$  is the dimensionality of the data. The second element ( $t_{\text{train}}$ ,  $t_{\text{valid}}$ ,  $t_{\text{test}}$ ) of each tuple is the corresponding  $N$ -dimensional vector of integers, containing the true class labels.

Here's a visualisation of the first 8 digits of the trainingset:

```
In [3]: def plot_digits(data, numcols, shape=(28,28)):
        numdigits = data.shape[0]
        numrows = int(numdigits/numcols)
        for i in range(numdigits):
            plt.subplot(numrows, numcols, i)
            plt.axis('off')
            plt.imshow(data[i].reshape(shape), interpolation='nearest', cmap='Greys')
        plt.show()

plot_digits(x_train[0:8], numcols=4)
```

```
/Library/Python/2.7/site-packages/matplotlib/axes/_subplots.py:69: MatplotlibDeprecationWarning: The use of 0 (which ends up being the _last_ sub-plot) is deprecated in 1.4 and will raise an error in 1.5
mplDeprecation)
```



In *multiclass* logistic regression, the conditional probability of class label  $j$  given the image  $\mathbf{x}$  for some datapoint is given by:

$$\log p(t = j \mid \mathbf{x}, \mathbf{b}, \mathbf{W}) = \log q_j - \log Z$$

where  $\log q_j = \mathbf{w}_j^T \mathbf{x} + b_j$  (the log of the unnormalized probability of the class  $j$ ), and  $Z = \sum_k q_k$  is the normalizing factor.  $\mathbf{w}_j$  is the  $j$ -th column of  $\mathbf{W}$  (a matrix of size  $784 \times 10$ ) corresponding to the class label,  $b_j$  is the  $j$ -th element of  $\mathbf{b}$ .

Given an input image, the multiclass logistic regression model first computes the intermediate vector  $\log \mathbf{q}$  (of size  $10 \times 1$ ), using  $\log q_j = \mathbf{w}_j^T \mathbf{x} + b_j$ , containing the unnormalized log-probabilities per class.

The unnormalized probabilities are then normalized by  $Z$  such that  $\sum_j p_j = \sum_j \exp(\log p_j) = 1$ . This is done by  $\log p_j = \log q_j - \log Z$  where  $Z = \sum_j \exp(\log q_j)$ . This is known as the *softmax* transformation, and is also used as a last layer of many classification neural network models, to ensure that the output of the network is a normalized distribution, regardless of the values of second-to-last layer ( $\log \mathbf{q}$ )

Warning: when computing  $\log Z$ , you are likely to encounter numerical problems. Save yourself countless hours of debugging and learn the [log-sum-exp trick](https://hips.seas.harvard.edu/blog/2013/01/09/computing-log-sum-exp/) (<https://hips.seas.harvard.edu/blog/2013/01/09/computing-log-sum-exp/>).

The network's output  $\log \mathbf{p}$  of size  $10 \times 1$  then contains the conditional log-probabilities  $\log p(t = j \mid \mathbf{x}, \mathbf{b}, \mathbf{W})$  for each digit class  $j$ . In summary, the computations are done in this order:

$$\mathbf{x} \rightarrow \log \mathbf{q} \rightarrow Z \rightarrow \log \mathbf{p}$$

Given some dataset with  $N$  independent, identically distributed datapoints, the log-likelihood is given by:

$$\mathcal{L}(\mathbf{b}, \mathbf{W}) = \sum_{n=1}^N \mathcal{L}^{(n)}$$

where we use  $\mathcal{L}^{(n)}$  to denote the partial log-likelihood evaluated over a single datapoint. It is important to see that the log-probability of the class label  $t^{(n)}$  given the image, is given by the  $t^{(n)}$ -th element of the network's output  $\log \mathbf{p}$ , denoted by  $\log p_{t^{(n)}}$ :

$$\mathcal{L}^{(n)} = \log p(t = t^{(n)} \mid \mathbf{x} = \mathbf{x}^{(n)}, \mathbf{b}, \mathbf{W}) = \log p_{t^{(n)}} = \log q_{t^{(n)}} - \log Z^{(n)}$$

where  $\mathbf{x}^{(n)}$  and  $t^{(n)}$  are the input (image) and class label (integer) of the  $n$ -th datapoint, and  $Z^{(n)}$  is the normalizing constant for the distribution over  $t^{(n)}$ .

## 1.1 Gradient-based stochastic optimization

### 1.1.1 Derive gradient equations (20 points)

Derive the equations for computing the (first) partial derivatives of the log-likelihood w.r.t. all the parameters, evaluated at a *single* datapoint  $n$ .

You should start deriving the equations for  $\frac{\partial \mathcal{L}^{(n)}}{\partial \log q_j}$  for each  $j$ . For clarity, we'll use the shorthand  $\delta_j^q = \frac{\partial \mathcal{L}^{(n)}}{\partial \log q_j}$ .

For  $j = t^{(n)}$ :  $\delta_{t^{(n)}}^q = \frac{\partial \mathcal{L}^{(n)}}{\partial \log p_{t^{(n)}}} \frac{\partial \log p_{t^{(n)}}}{\partial \log q_{t^{(n)}}}$

- $\frac{\partial \mathcal{L}^{(n)}}{\partial \log Z} \frac{\partial \log Z}{\partial \log q_j} \frac{\partial \log q_j}{\partial \log p_{t^{(n)}}} \frac{\partial \log p_{t^{(n)}}}{\partial \log q_{t^{(n)}}} = 1$   
 $\cdot 1 - \frac{\partial \log Z}{\partial \log q_j} \frac{\partial \log q_j}{\partial \log p_{t^{(n)}}} \frac{\partial \log p_{t^{(n)}}}{\partial \log q_{t^{(n)}}} = 1 - \frac{\partial \log Z}{\partial \log q_j} \frac{\partial \log q_j}{\partial \log p_{t^{(n)}}}$

For  $j \neq t^{(n)}$ :  $\delta_j^q = \frac{\partial \mathcal{L}^{(n)}}{\partial \log Z} \frac{\partial \log Z}{\partial \log q_j} \frac{\partial \log q_j}{\partial \log p_{t^{(n)}}} \frac{\partial \log p_{t^{(n)}}}{\partial \log q_{t^{(n)}}} = - \frac{\partial \log Z}{\partial \log q_j} \frac{\partial \log q_j}{\partial \log p_{t^{(n)}}}$

Complete the above derivations for  $\delta_j^q$  by furtherly developing  $\frac{\partial \log Z}{\partial \log q_j}$  and  $\frac{\partial \log q_j}{\partial \log p_{t^{(n)}}}$ . Both are quite simple. For these it doesn't matter whether  $j = t^{(n)}$  or not.

Given your equations for computing the gradients  $\delta_j^q$  it should be quite straightforward to derive the equations for the gradients of the parameters of the model,  $\frac{\partial \mathcal{L}^{(n)}}{\partial W_{ij}}$  and  $\frac{\partial \mathcal{L}^{(n)}}{\partial b_j}$ . The gradients for the biases  $\mathbf{b}$  are given by:

$$\frac{\partial \mathcal{L}^{(n)}}{\partial b_j} = \frac{\partial \mathcal{L}^{(n)}}{\partial \log q_j} \frac{\partial \log q_j}{\partial b_j} = \delta_j^q \cdot 1 = \delta_j^q$$

The equation above gives the derivative of  $\mathcal{L}^{(n)}$  w.r.t. a single element of  $\mathbf{b}$ , so the vector  $\nabla_{\mathbf{b}} \mathcal{L}^{(n)}$  with all derivatives of  $\mathcal{L}^{(n)}$  w.r.t. the bias parameters  $\mathbf{b}$  is:

$$\nabla_{\mathbf{b}} \mathcal{L}^{(n)} = \delta^q$$

where  $\delta^q$  denotes the vector of size  $10 \times 1$  with elements  $\delta_j^q$ .

The (not fully developed) equation for computing the derivative of  $\mathcal{L}^{(n)}$  w.r.t. a single element  $W_{ij}$  of  $\mathbf{W}$  is:

$$\frac{\partial \mathcal{L}^{(n)}}{\partial W_{ij}} = \frac{\partial \mathcal{L}^{(n)}}{\partial \log q_j} \frac{\partial \log q_j}{\partial W_{ij}} = \delta_j^q \frac{\partial \log q_j}{\partial W_{ij}}$$

What is  $\frac{\partial \log q_j}{\partial W_{ij}}$ ? Complete the equation above.

If you want, you can give the resulting equation in vector format ( $\nabla_{\mathbf{W}_j} \mathcal{L}^{(n)} = \dots$ ), like we did for  $\nabla_{\mathbf{b}} \mathcal{L}^{(n)}$ .

**Answer:**

[insert answer in this cell]

### 1.1.2 Implement gradient computations (10 points)

Implement the gradient calculations you derived in the previous question. Write a function `logreg_gradient(x, t, w, b)` that returns the gradients  $\nabla_{\mathbf{w}_j} \mathcal{L}^{(n)}$  (for each  $j$ ) and  $\nabla_{\mathbf{b}} \mathcal{L}^{(n)}$ , i.e. the first partial derivatives of the log-likelihood w.r.t. the parameters  $\mathbf{W}$  and  $\mathbf{b}$ , evaluated at a single datapoint  $(\mathbf{x}, t)$ . The computation will contain roughly the following intermediate variables:

$$\log \mathbf{q} \rightarrow Z \rightarrow \log \mathbf{p}, \delta^q$$

followed by computation of the gradient vectors  $\nabla_{\mathbf{w}_j} \mathcal{L}^{(n)}$  (contained in a  $784 \times 10$  matrix) and  $\nabla_{\mathbf{b}} \mathcal{L}^{(n)}$  (a  $10 \times 1$  vector).

In [ ]:

### 1.1.3 Stochastic gradient descent (10 points)

Write a function `sgd_iter(x_train, t_train, w, b)` that performs one iteration of stochastic gradient descent (SGD), and returns the new weights. It should go through the trainingset once in randomized order, call `logreg_gradient(x, t, w, b)` for each datapoint to get the gradients, and update the parameters using a small learning rate (e.g.  $1\text{E-}4$ ). Note that in this case we're maximizing the likelihood function, so we should actually performing gradient **ascent**... For more information about SGD, see Bishop 5.2.4 or an online source (i.e. [https://en.wikipedia.org/wiki/Stochastic\\_gradient\\_descent](https://en.wikipedia.org/wiki/Stochastic_gradient_descent) ([https://en.wikipedia.org/wiki/Stochastic\\_gradient\\_descent](https://en.wikipedia.org/wiki/Stochastic_gradient_descent))).

In [ ]:

## 1.2. Train

### 1.2.1 Train (10 points)

Perform a handful of training iterations through the trainingset. Plot (in one graph) the conditional log-probability of the trainingset and validation set after each iteration.

In [ ]:

### 1.2.2 Visualize weights (10 points)

Visualize the resulting parameters  $\mathbf{W}$  after a few iterations through the training set, by treating each column of  $\mathbf{W}$  as an image. If you want, you can use or edit the `plot_digits(...)` above.

In [ ]:

### 1.2.3. Visualize the 8 hardest and 8 easiest digits (10 points)

Visualize the 8 digits in the validation set with the highest probability of the true class label under the model. Also plot the 8 digits that were assigned the lowest probability. Ask yourself if these results make sense.

In [ ]:

## Part 2. Multilayer perceptron

You discover that the predictions by the logistic regression classifier are not good enough for your application: the model is too simple. You want to increase the accuracy of your predictions by using a better model. For this purpose, you're going to use a multilayer perceptron (MLP), a simple kind of neural network. The perceptron will have a single hidden layer  $\mathbf{h}$  with  $L$  elements. The parameters of the model are  $\mathbf{V}$  (connections between input  $\mathbf{x}$  and hidden layer  $\mathbf{h}$ ),  $\mathbf{a}$  (the biases/intercepts of  $\mathbf{h}$ ),  $\mathbf{W}$  (connections between  $\mathbf{h}$  and  $\log q$ ) and  $\mathbf{b}$  (the biases/intercepts of  $\log q$ ).

The conditional probability of the class label  $j$  is given by:

$$\log p(t = j \mid \mathbf{x}, \mathbf{b}, \mathbf{W}) = \log q_j - \log Z$$

where  $q_j$  are again the unnormalized probabilities per class, and  $Z = \sum_j q_j$  is again the probability normalizing factor. Each  $q_j$  is computed using:

$$\log q_j = \mathbf{w}_j^T \mathbf{h} + b_j$$

where  $\mathbf{h}$  is a  $L \times 1$  vector with the hidden layer activations (of a hidden layer with size  $L$ ), and  $\mathbf{w}_j$  is the  $j$ -th column of  $\mathbf{W}$  (a  $L \times 10$  matrix). Each element of the hidden layer is computed from the input vector  $\mathbf{x}$  using:

$$h_j = \sigma(\mathbf{v}_j^T \mathbf{x} + a_j)$$

where  $\mathbf{v}_j$  is the  $j$ -th column of  $\mathbf{V}$  (a  $784 \times L$  matrix),  $a_j$  is the  $j$ -th element of  $\mathbf{a}$ , and  $\sigma(\cdot)$  is the so-called sigmoid activation function, defined by:

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

Note that this model is almost equal to the multiclass logistic regression model, but with an extra 'hidden layer'  $\mathbf{h}$ . The activations of this hidden layer can be viewed as features computed from the input, where the feature transformation ( $\mathbf{V}$  and  $\mathbf{a}$ ) is learned.

### 2.1 Derive gradient equations (20 points)

State (shortly) why  $\nabla_{\mathbf{b}} \mathcal{L}^{(n)}$  is equal to the earlier (multiclass logistic regression) case, and why  $\nabla_{\mathbf{w}_j} \mathcal{L}^{(n)}$  is almost equal to the earlier case.

Like in multiclass logistic regression, you should use intermediate variables  $\delta_j^q$ . In addition, you should use intermediate variables  $\delta_j^h = \frac{\partial \mathcal{L}^{(n)}}{\partial h_j}$ .

Given an input image, roughly the following intermediate variables should be computed:

$$\log \mathbf{q} \rightarrow Z \rightarrow \log \mathbf{p} \rightarrow \delta^q \rightarrow \delta^h$$

$$\text{where } \delta_j^h = \frac{\partial \mathcal{L}^{(n)}}{\partial h_j}.$$

Give the equations for computing  $\delta^h$ , and for computing the derivatives of  $\mathcal{L}^{(n)}$  w.r.t.  $\mathbf{W}$ ,  $\mathbf{b}$ ,  $\mathbf{V}$  and  $\mathbf{a}$ .

You can use the convenient fact that  $\frac{\partial}{\partial x} \sigma(x) = \sigma(x)(1 - \sigma(x))$ .

**Answer:**

[insert answer in this Markdown cell]

## 2.2 MAP optimization (10 points)

You derived equations for finding the *maximum likelihood* solution of the parameters. Explain, in a few sentences, how you could extend this approach so that it optimizes towards a *maximum a posteriori* (MAP) solution of the parameters, with a Gaussian prior on the parameters.

**Answer:**

[insert answer in this Markdown cell]

## 2.3. Implement and train a MLP (15 points)

Implement a MLP model with a single hidden layer, and code to train the model.

In [ ]:

### 2.3.1. Less than 250 misclassifications on the test set (10 bonus points)

You receive an additional 10 bonus points if you manage to train a model with very high accuracy: at most 2.5% misclassified digits on the test set. Note that the test set contains 10000 digits, so your model should misclassify at most 250 digits. This should be achievable with a MLP model with one hidden layer. See results of various models at : <http://yann.lecun.com/exdb/mnist/index.html>. To reach such a low accuracy, you probably need to have a very high  $L$  (many hidden units), probably  $L > 200$ , and apply a strong Gaussian prior on the weights. In this case you are allowed to use the validation set for training. You are allowed to add additional layers, and use convolutional networks, although that is probably not required to reach 2.5% misclassifications.

In [ ]: