



Bright



Dark



Blues



Grays



Night

Project 3: Softmax Regression and NN

In this project, we will explore softmax regression and neural networks. The project files are available for download in the p3 folder on the course Files page.

Files to turn in:

softmax.py	Implementation of softmax regression
nn.py	Implementation of fully connected neural network
partners.txt	Lists the full names of all members in your team.
writeup.pdf	Answers all the written questions in this assignment
Files you created for Extra-Credits	

You will be using helper functions in the following .py files and datasets:

utils.py	Helper functions for Softmax Regression and NN
data/*	Training and dev data for SR and NN

Please **do not** change the file names listed above. Submit partners.txt even if you do the project alone. **Only one member** in a team (even if it is a cross-section team) is supposed to submit the project.

AutogradingPlease **do not** change the names of any provided functions or classes within the code, or you will wreak havoc on the autograder. However, the correctness of your implementation -- not the autograder's output -- will be the final judge of your score. If necessary, we will review and grade assignments individually to ensure that you receive due credit for your work. Please **do not** import (potentially unsafe) system-related modules such as sys, exec, eval... Otherwise the autograder will assign a zero point without grading.

Part I Softmax Regression [50%]

Files to edit/turn in for this part

softmax.py
writeup.pdf

The goal of this part of the project is to implement Softmax Regression in order to classify the MNIST digit dataset. Softmax Regression is essentially a two-layer neural network where the output layer applies the Softmax cost function, a multiclass generalization of the logistic cost function.

In logistic regression, we have a hypothesis function of the form

$$P[y = 1] = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}}$$

where \vec{w} is our weight vector. Like the hyperbolic tangent function, the logistic function is also a sigmoid function with the characteristic 's'-like shape, though it has a range of (0, 1) instead of (-1, 1). Note that this is technically not a classifier since it returns probabilities instead of a predicted class, but it's easy to turn it into a classifier by simply choosing the class with the highest probability.

Since logistic regression is used for binary classification, it is easy to see that:

$$\begin{aligned} P[y = 1] &= \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}} \\ &= \frac{e^{\vec{w} \cdot \vec{x}}}{e^{\vec{w} \cdot \vec{x}} + 1} \\ &= \frac{e^{\vec{w} \cdot \vec{x}}}{e^{\vec{w} \cdot \vec{x}} + e^{\vec{0} \cdot \vec{x}}} \end{aligned}$$

Similarly,

$$\begin{aligned} P[y = 0] &= 1 - \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}} \\ &= \frac{e^{\vec{w} \cdot \vec{x}} + 1}{e^{\vec{w} \cdot \vec{x}} + 1} - \frac{e^{\vec{w} \cdot \vec{x}}}{e^{\vec{w} \cdot \vec{x}} + 1} \\ &= \frac{e^{\vec{0} \cdot \vec{x}}}{e^{\vec{w} \cdot \vec{x}} + e^{\vec{0} \cdot \vec{x}}} \end{aligned}$$

From this form it appears that we can assign the vector $\vec{w}_1 = \vec{w}$ as the weight vector for class 1 and $\vec{w}_0 = \vec{0}$ as the weight vector for class 0. Our probability formulas are now unified into one equation:

$$P[y = i] = \frac{e^{\vec{w}_i \cdot \vec{x}}}{\sum_j e^{\vec{w}_j \cdot \vec{x}}}$$

This immediately motivates generalization to classification with more than 2 classes. By assigning a separate weight vector \vec{w}_i to each class, for each example \vec{x} we can predict the probability that it is class i , and again we can classify by choosing the most probable class. A more compact way of representing the values $\vec{w}_i \cdot \vec{x}$ is $W\vec{x}$ where each row i of W is \vec{w}_i . We can also represent a dataset $\{\vec{x}_i\}$ with a matrix X where each column is a single example.

Qsr1 (10%)

- (1) Show that the probabilities sum to 1.
- (2) What are the dimensions of W ? X ? WX ?

We can also train on this model with an appropriate loss function. The Softmax loss function is given by

$$L(W) = - \left[\sum_{i=1}^m \sum_{k=1}^K \mathbf{1}\{y_i = k\} \log \frac{e^{\vec{w}_k \cdot \vec{x}_i}}{\sum_{j=1}^K e^{\vec{w}_j \cdot \vec{x}_i}} \right]$$

where m is the number of examples, k is the number of classes, and $\mathbf{1}\{y_i = k\}$ is an indicator variable that equals 1 when the statement inside the brackets is true, and 0 otherwise. The gradient (which you will not derive) is given by:

$$\nabla_{\vec{w}_k} L(W) = - \sum_{i=1}^m \left[\vec{x}_i (\mathbf{1}\{y_i = k\} - P[y_i = k]) \right]$$

Note that the indicator and the probabilities can be represented as matrices, which makes the code for the loss and the gradient very simple. (See [here](#) for more details)

softmax.py contains a mostly-complete implementation of Softmax Regression. A code stub also has been provided in run_softmax.py. Once you correctly implement the incomplete portions of softmax.py, you will be able to run run_softmax.py in order to classify the MNIST digits.

Qsr2 (15%)

- (1) Complete the implementation of the cost function.
- (2) Complete the implementation of the predict function.

Check your implementation by running:

```
>>> python run_softmax.py
```

The output should be:

RUNNING THE L-BFGS-B CODE

* * *

Machine precision = 2.220D-16

N = 7840 M = 10

At X0 0 variables are exactly at the bounds

At iterate 0 f= 2.30259D+00 |proj g|= 6.37317D-02

At iterate 1 f= 1.52910D+00 |proj g|= 6.91122D-02

At iterate 2 f= 7.72038D-01 |proj g|= 4.43378D-02

...

At iterate 401 f= 2.19686D-01 |proj g|= 2.52336D-04

At iterate 402 f= 2.19665D-01 |proj g|= 2.04576D-04

* * *

Tit = total number of iterations

Tnf = total number of function evaluations

Tnint = total number of segments explored during Cauchy searches

Skip = number of BFGS updates skipped

Nact = number of active bounds at final generalized Cauchy point

Projg = norm of the final projected gradient

F = final function value

* * *

N Tit Tnf Tnint Skip Nact Projg F

7840 402 431 1 0 0 2.046D-04 2.197D-01

F = 0.21966482316858085

STOP: TOTAL NO. of ITERATIONS EXCEEDS LIMIT

Cauchy time 0.000E+00 seconds.

Subspace minimization time 0.000E+00 seconds.

Line search time 0.000E+00 seconds.

Total User time 0.000E+00 seconds.

Accuracy: 93.99%

Qsr3 (15%)

In the cost function, we see the line

$W_X = W_X - \text{np.max}(W_X)$

This means that each entry is reduced by the largest entry in the matrix.

(1) Show that this does not affect the predicted probabilities.

(2) Why might this be an optimization over using W_X ? Justify your answer.

Qsr4 (10%)

Use the learningCurve function in runClassifier.py to plot the accuracy of the classifier as a function of the number of examples seen. Include the plot in your write-up. Do you observe any overfitting or underfitting? Discuss and explain what you observe.

Part II NN [50% and Extra 15%]

Files to edit/turn in.

nn.py

writeup.pdf

files created for the Q2 extra credits

In this part of the project, we'll implement a fully-connected neural network in general for the MNIST dataset. For **Qnn1** you will complete nn.py. For **Qnn2**, create your own files.

A code stub also has been provided in run_nn.py. Once you correctly implement the incomplete portions of nn.py, you will be able to run run_nn.py in order to classify the MNIST digits.

The dataset is included under ./data/. You will be using the following helper functions in "utils.py". In the following description, let K , N , d denote the number of classes, number of samples and number of features.

- Selected functions in "utils.py"

```
loadMNIST(image_file, label_file) #returns data matrix X with shape (d, N) and labels
onehot(labels) # encodes labels into one hot style with shape (K,N)
acc(pred_label, Y) # calculate the accuracy of prediction given ground truth Y, w
data_loader(X, Y=None, batch_size=64, shuffle=False) # Iterator that yield X,Y wi
```

Qnn1 (45% for Qnn1.1, 1.2, 1.3 (15% each) and 5% for Qnn 1.4) Implement the NN

The scaffold has been built for you. Initialize the model and print the architecture with the following:

```
>>> from nn import NN, Relu, Linear, SquaredLoss
>>> from utils import data_loader, acc, save_plot, loadMNIST, onehot
>>> model = NN(Relu(), SquaredLoss(), hidden_layers=[128,128])
>>> model.print_model()
```

Two activation functions (Relu, Linear) and self.predict(X) have been implemented for you.

Qnn1.1 Implement squared loss cost functions (TODO 0 & TODO 1)

Assume \bar{Y} is the output of the last layer before loss calculation (without activation), which is a K-by-N matrix. Y is the one hot encoded ground truth of the same shape. Implement the following loss function and its gradient (You need to calculate and implement the gradient of the loss function yourself) (Notice that the loss functions are normalized by batch_size N):

$$L(Y, \bar{Y}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \|\bar{Y}_i - Y_i\|^2 = \frac{1}{2N} \|\bar{Y} - Y\|_{fro}^2, \text{ where } Y_i \text{ is the } i\text{-th column of } Y.$$

Typically we would use cross entropy loss, the formula of which is provided for your reference (but you are only required to implement for squared loss):

$$L_{CE}(Y, \bar{Y}) = \frac{1}{N} \sum_{i=1}^N NLL(Y_i, \bar{Y}_i), \text{ where}$$

$$NLL(y, \bar{y}) = -\log\left(\sum_{j=1}^K y_j \frac{\exp(\bar{y}_j)}{\sum_{k=1}^K \exp(\bar{y}_k)}\right).$$

Qnn1.2. Compute the gradients (TODO 2 & TODO 3)

Implement the forward pass (TODO 2) and back propagation (TODO 3) for gradient calculation. Use "activation.activate" and "activation.backprop_grad" in your code so that your gradient computation works for different choices of activation functions.

Do the following to see if the loss goes down.

```
>>> x_train, label_train = loadMNIST('data/train-images.idx3-ubyte', 'data/train-
>>> x_test, label_test = loadMNIST('data/t10k-images.idx3-ubyte', 'data/t10k-labe
```

```

>>> y_train = onehot(label_train)
>>> y_test = onehot(label_test)

>>> model = NN(Relu(), SquaredLoss(), hidden_layers=[128, 128], input_d=784, outp
>>> model.print_model()
>>> training_data, dev_data = {"X":x_train, "Y":y_train}, {"X":x_test, "Y":y_test
>>> from run_nn import train_1pass
>>> model, plot_dict = train_1pass(model, training_data, dev_data, learning_rate=

```

Qnn1.3 Run in epochs

An epoch is a full pass of the training data. Run `run_nn.py`.

```
>>> python3 run_nn.py
```

Report your final accuracy on the dev set. You can either use the default setting or tune the architecture (number of layers, size of layers and loss function) and hyperparameters (lr, batch_size, max_epoch).

Qnn1.4 (No implementation needed for this question). When initializing the weight matrix, in some cases it may be appropriate to initialize the entries as small random numbers rather than all zeros. Give one reason why this may be a good idea.

Qnn2 (Extra-Credit 15%) Try something new.

Choose one of the following directions (outside research may be required) for further exploration (Feel free to copy `nn.py`, `utils.py` and `run_nn.py` as a starting point. **Make sure** that your code for Qnn2 is separated from your code for Qnn1):

(1) Do dimension reduction with PCA. Try with different dimensions. Can you observe the trade-off in time and acc? Plot training time v.s. dimension, testing time v.s dimension and acc v.s. dimension. Visualize the principal components.

(2) Improve your results with ensemble methods. Describe your implementation. Can you observe improved performance compared with that of Q4? Why?
(http://ciml.info/dl/v0_99/ciml-v0_99-ch13.pdf)

(3) Implement a new optimizer (By implementing a different `self.update` for the NN class). Compare with the original SGD optimizer. You can read about the optimizers in (<http://ruder.io/deep-learning-optimization-2017/>). Does this new method take less number of

samples to converge? Does this new method take less time to converge?

Qnn2.1 Explain what you did and what you found. Comment the code so that it is easy to follow. Support your results with plots and numbers. Provide the implementation so we can replicate **your results**.