# KTH Royal Institute of Technology

## DD2424
### Deep Learning in Data Science

### Year 2017 - 2018

# Assignment 1

Picó Oristrell Sandra                    940531 − 2308    sandrapo@kth.se

*Professors :*
Sullivan Josephine

April 11, 2018

# Contents

The main intention of this assignment is to train and test a one layer network to classify images from the CIFAR-10 dataset. The algorithm used is a mini-batch gradient descent with a cost function that computes the cross-entropy loss and with an L$_2$ regularization term on the weight matrix.

In order to explain and show the results obtained in a proper way I will divide this report into different subsections. The first section will be focused on the background of the assignment where I will explain the main functions of the classifier and the data that we are analyzing. To continue, I will specified the functions implemented. As I am doing the assignment in Matlab, the functions have the same structure as in the assignment description.All of them are attached in the appendix section. Finally I will report all the results obtained in the exercise 1.

# 1 Background

As mentioned above, in this section I will explain a little bit about the background needed for this assignment.

## 1.1 Mathematical background:

The network that we are going to implement will be a classification problem where the output will be a vector of probabilities.

Figure 1: Main structure of the network used

Then, the predicted class will correspond to the label with higher probability, i.e the *arg max* of the probability's vector. Nevertheless, in order to compute the output vector we will use the following formulas:

$$s = Wx + b \tag{1}$$

$$\mathbf{p} = \frac{exp(s)}{1^T exp(s)} \tag{2}$$

This means that during the training phase we will learn the parameters $W$ and $b$. To set this parameters we will solve the following optimization problem:

$$W^*, \mathbf{b}^* = argmin_{W,b} J(D, \lambda, W, b) \tag{3}$$

Where $J$ is a cost function that minimize the cross-entropy plus a regularization term on $W$.

$$J(D, \lambda, W, b) = \frac{1}{|D|} \sum_{x,y \in D} \log(y^T p) + \lambda \sum_{i,j} W_{ij}^2 \tag{4}$$

The problem will be solve via mini-batch gradient descent algorithm.

## 1.2 Data structure:

As mentioned before, the data that we are going to use is from the CIFAR-10 dataset. The dataset has a size of 10000 x 3072 following the next structure:
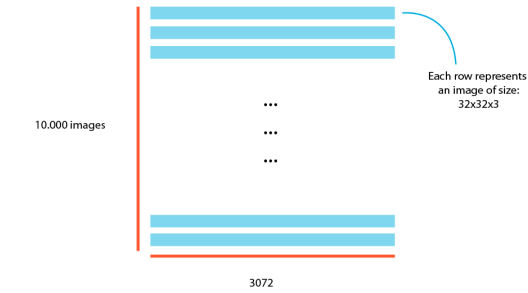


Figure 2: Data structure of CIFAR-10 dataset

# 2 Implementation of the functions

To be able to train the multi-linear classifier we will need to write some sub-functions. The purpose of this section is to explain them. Nevertheless, you can find the code of each attached in the appendix section.

- **function [X,Y,y] = LoadBatch(filename)**

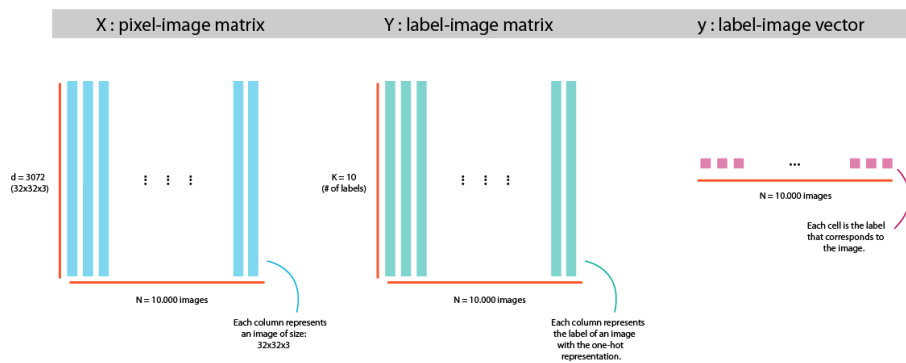  This function reads a filename and returns the image and the label data in the following format:



Figure 3: X,Y and y structure

- **function P = EvaluateClassifier(X,W,b)**

  This function generates the P matrix. As explained above, the P matrix will be a matrix where each column will show the probability for each label for a specific image.

- **function J = ComputeCost(X,Y,W,b,lambda)**

  This function computes the cost based on the cross-entropy loss and the regularization term on $W$. This function will be called when we implement the gradient descent algorithm.

- **function acc = ComputeAccuracy(P,y)**

  To be able to compute the accuracy of our classifier we need to count how many outputs are correct , comparing with the truth labels and compute the percentage.

- **function [grad_W,grad_b] = ComputeGradients(X,Y,P,W,lambda)**

In this function we are going to compute the analytic gradients. The calculus are based on the lecture's slides. This gradients are first checked through the numerical gradients computed in the ComputeGradsNumSlow.

- **function [Wstar,bstar] = MiniBatchGD(X,Y,GDparams,W,b,lambda)**

  Finally, this function will implement the vanilla version of the mini-batch gradient descent algorithm. This algorithm will learn the network parameters: W and b. The main structure of the code is specified in the description of the assignment.

  GDparams is an object that contains the parameter values needed for the MiniBatchGD function: n_batch,eta and n_epochs.

# 3 Exercise 1:

## 3.1 Check the gradients

The first top exercise is to be sure that the analytical gradients are properly calculated. In order to do so, I decide to compute the relative error specified in the description. This error is computed as follows:

$$\frac{\mid g_a - g_n \mid}{max(eps, \mid g_a \mid + \mid g_n \mid)} \tag{5}$$

where $g_n$ is the numerically computed gradient and $g_a$ is the analytically computed gradient. If the relative error obtained is $<$ 1e-6 then we are sure that the implementation is correct. The implementation can be found in the appendix section (Assignment1_1).

## 3.2 Loss, accuracy and weights visualization

Once I was sure that the gradients were correctly implemented I plot the accuracy, the loss and the weights visualization per each configuration specified in the assignment description. The implementation can be found in the appendix section (Assignment1_2). The obtained results are the following ones:

- **No-Regularization:**

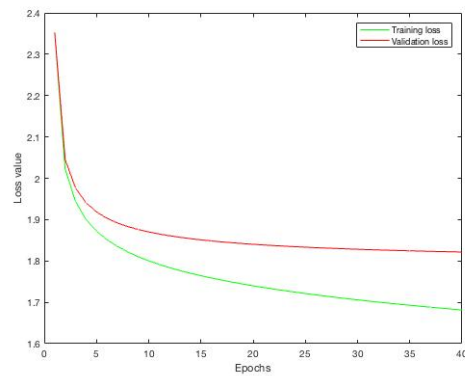  - **lambda=0, n epochs=40, n batch=100, eta=.1**
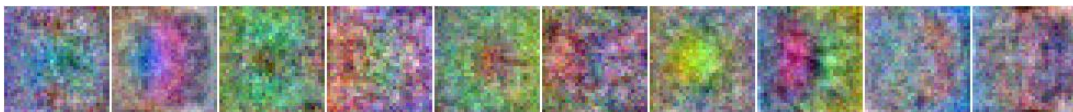
    * *Loss:*

* *Weights:*



* *Accuracy:* Training Accuracy:31.59% Test Accuracy:26.2%

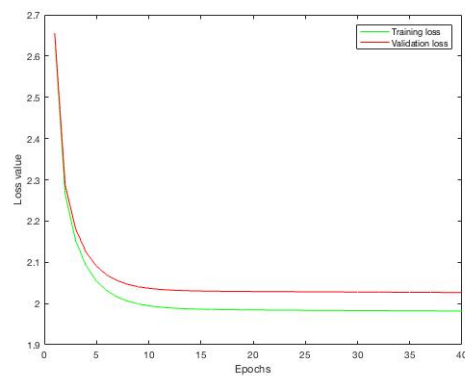– **lambda=0, n epochs=40, n batch=100, eta=.01**

* *Loss:*



* *Weights:*



* *Accuracy:* Training Accuracy:41.56% Test Accuracy:36.87%

• **Regularization:**

– **lambda=.1, n epochs=40, n batch=100, eta=.01**

* *Loss:*



* *Weights:*

* *Accuracy:* Training Accuracy:34.18% Test Accuracy:33.37%

– **lambda=1, n epochs=40, n batch=100, eta=.01**

* *Loss:*



* *Weights:*



* *Accuracy:* Training Accuracy:22.24% Test Accuracy:21.92%

## 3.3 Conclusions

Regarding the results obtained in the previous section we can state two main conclusions regarding the learning rate and the regularization term. A high learning rate implies that the loss cost function specified will converge faster. However, that also means that it will be unstable. This behavior can be found in the first loss graph. Better performance is achieved with a low learning rate.

On the other hand, regularization term implies reducing the over-fitting effect. In the two last cases we have seen that the space between validation loss and training loss is less than in the two cases where we did not have regularization. As we don't have over-fitting,we have better generalization.

# Appendices

## A Functions implementation:

### A.1 LoadBatch function

```
1  function[X,Y,y] = LoadBatch(filename)
2      %Open the dataset
3      addpath Datasets/;
4      dataset = load(filename);
5      %X contains the image pixel data. Size 3072x10000. Double type between
6      %0 and 1.
7      %Reshape the pixel data into a 3072x10000 array
8      I = reshape(dataset.data.',32,32,3,10000);
9      I = permute(I,[2,1,3,4]);
10     X = reshape(I,[3072,10000]);
11     %Normalize the vector
12     X = double(X);
13     X = X./255.0;
14
15     %y is a vector of 10000 containing the label for each image.
16     %better to encode it between 1-10 instead of 0-9.
17     y = dataset.labels;
18     for i = 1:10000
19         y(i,1) = y(i,1) + 1;
20     end
21
22     %Y is a 10x10000 array with the one_hot representation.
23     %Extract the labels dataset
24     Y = dataset.labels;
25     %One-hot representation
26     num_labels = 10;      %Number of labels
27     Y_one_hot = zeros(size(Y,1),num_labels);
28     Y_one_hot = reshape(Y_one_hot,[10,10000]); %Reshape the one_hot representation.
29     for i = 1:10000
30         row = Y_one_hot(:,i);
31         label = Y(i,:);
32         row(label+1) = 1;
33         Y_one_hot(:,i) = row;
34     end
35     Y = Y_one_hot;
36
37  return
```

### A.2 EvaluateClassifier function

```
1  function P = EvaluateClassifier(X,W,b)
2      %Be sure that it is double.
3      W = double(W);
4      X = double(X);
5      b = double(b);
6      %Softmax to generate a vector with the probability of each label.
7      s = W*X + b;
8      K = size(b);
9      K = K(1);
10     %Compute the P matrix.
11     %Each column of P represents the probability of each label /image.
```

```matlab
12      P = exp(s)./((ones(1,K,'double'))*exp(s));
13  return
```

## A.3   ComputeCost function

```matlab
1   function J = ComputeCost(X,Y,W,b,lambda)
2       %Size values
3       K = 10;
4       %In this case, N will represent the D subset.
5       [d,N] = size(X);
6
7       %Compute P
8       W = double(W);
9       X = double(X);
10      b = double(b);
11      s = W*X + b;
12      K = size(b);
13      K = K(1);
14      P = exp(s)./((ones(1,K,'double'))*exp(s));
15
16      %Cross-Entropy loss
17      crossEntropy_term = 0;
18      for i = 1:N
19          y = Y(:,i);
20          p = P(:,i);
21          py = y.'*p;
22          crossEntropy_term = crossEntropy_term -log(py);
23      end
24
25      %Regularization term.
26      reg_term = 0;
27      for row = 1:K
28          for column = 1:d
29              reg_term = reg_term + W(row,column)^2;
30          end
31      end
32      reg_term = lambda*reg_term;
33
34      %Calculate J.
35      J = 1./abs(N)*crossEntropy_term + reg_term;
36
37  return
```

## A.4   ComputeAccuracy function

```matlab
1   function acc = ComputeAccuracy(P,y)
2       %size P = KxN
3       [K,N] = size(P);
4       %size y = Nx1
5       %To store the prediction using the P matrix.
6       pred = zeros(N,1,'double');
7       %Compute the result using the P matrix.
8       for i = 1: N
9           [argvalue, argmax] = max(P(:,i));
10          pred(i,1) = argmax;
11      end
12      correct = 0;
13      for i = 1: N
```

```
14          if pred(i,1) == y(i,1)
15              correct = correct + 1;
16          end
17      end
18      acc = (correct/N)*100;
19  return
```

## A.5   ComputeGradients function

```
1  function [grad_W,grad_b] = ComputeGradients(X,Y,P,W,lambda)
2      % grad_W: gradient matrix of the cost J respect to W.
3      % grad_b: gradient vector of the cost J respect to b.
4
5      %Initialize the gradients with the respective sizes.
6      grad_W = zeros(size(W));
7      grad_b = zeros(size(W, 1), 1);
8
9      %Compute gradients from the slides (lectures).
10     for i = 1 : size(X, 2)
11         Xi = X(:, i);
12         Pi = P(:, i);
13         Yi = Y(:, i);
14         g = -Yi'/(Yi'*Pi)*(diag(Pi) - Pi*Pi');
15         grad_b = grad_b + g';
16         grad_W = grad_W + g'*Xi';
17     end
18
19     %Normalize the gradients.
20     grad_b = grad_b/size(X, 2);
21     grad_W = grad_W/size(X, 2) + (2*lambda)*W;
22  end
```

## A.6   MiniBatchGD function

```
1  function [Wstar,bstar] = MiniBatchGD(X,Y,GDparams,W,b,lambda)
2      %Size of X is : dxN
3      N = size(X, 2);
4
5      %GDparams is an object containing the parameter values n_batch, eta and
6      %n_epochs.
7      eta = GDparams.eta;
8      n_batch = GDparams.n_batch;
9      GDparams.n_epochs;
10
11     %Structure of the assignment description.
12     for j=1:N/n_batch
13         j_start = (j-1)*n_batch + 1;
14         j_end = j*n_batch;
15         inds = j_start:j_end;
16         Xbatch = X(:, inds);
17         Ybatch = Y(:, inds);
18         %EvaluateClassifier to be able to compute the gradients at each
19         %mini-batch computation.
20         P = EvaluateClassifier(Xbatch, W, b);
21         [grad_W, grad_b] = ComputeGradients(Xbatch, Ybatch, P, W, lambda);
22         %Update the values based with the learning-rate.
23         W = W - eta*grad_W;
24         b = b - eta*grad_b;
```

```
25        end
26        %Parameters that we would like to return
27        Wstar = W;
28        bstar = b;
29    end
```

# B  Exercise 1 implementation:

## B.1  Assignment 1.1:

```
1   function Assignment1_1
2       %%Exercise 1.1
3       %Extract data from the datafiles.
4       addpath Datasets/;
5       training_name = 'data_batch_1.mat';
6       test_name = 'test_batch.mat';
7       validation_name = 'data_batch_2.mat';
8       %X,Y and y data for each file.
9       [X_train,Y_train,y_train] = LoadBatch(training_name);
10      [X_test,Y_test,y_test] = LoadBatch(test_name);
11      [X_validation,Y_validation,y_validation] = LoadBatch(validation_name);
12
13      %Initialize the parameters of the model W and b. [size(W) = 10x3072
14      %size(b) = 10x1]
15      %Initialize each entry to have Gaussian random values with 0 mean and
16      %standard deviation 0.01.
17      c = 0;
18      a = 0.01;
19      K = size(Y_train,1);
20      d = size(X_train,1);
21      b = a.*randn(K,1)+c;
22      W = a.*randn(K,d)+c;
23
24      %Parameters to check:
25      lambda = 0;
26      batch = 100;
27      %Small positive number
28      eps = 1e-10;
29
30      %Compute the numerically gradients.
31      [grad_b_n, grad_W_n] = ComputeGradsNumSlow(X_train(:, 1 : batch), Y_train(:, 1 : ...
            batch), W, b, lambda, 1e-6);
32      %Generate P to be able to compute the gradients with our function
33      %ComputeGradients.
34      P = EvaluateClassifier(X_train(:, 1 : batch),  W, b);
35      %Analytical gradients.
36      [grad_W, grad_b] = ComputeGradients(X_train(:, 1 : batch), Y_train(:, 1 : batch), ...
            P,  W, lambda);
37      %Compare them to be sure that the gradients are computed properly.
38      %Compare them through the computation of the relative error.
39      %|ga-gn| / (max(eps,|ga|+|gn|)) where eps is a small positive number.
40      %Check this is small.
41      gradient_W = max(max(abs(grad_W_n - grad_W)./max(eps, abs(grad_W_n) + abs(grad_W))));
42      gradient_b = max(abs(grad_b_n - grad_b)./max(eps, abs(grad_b_n) + abs(grad_b)));
43
44      %Check if gradients are computed correctly.
45      %Small number choosed: 1e-6 (specified in the description of the
46      %assignment)
47
48      %Gradient W:
```

```
49        if gradient_W < 1e-6
50            fprintf("Correct gradient W!");
51        else
52            fprintf("Incorrect gradient W!");
53        end
54        %Gradient b:
55        if gradient_b < 1e-6
56            fprintf("Correct gradient b!");
57        else
58            fprintf("Incorrect gradient b!");
59        end
60    end
```

## B.2  Assignment 1.2:

```
1   function Assignment1_2
2       %Now we know for sure that the gradients calculations are ok.
3       %Exercise 1.2 : Training a multi-linear classifier.
4       addpath Datasets/;
5       training_name = 'data_batch_1.mat';
6       test_name = 'test_batch.mat';
7       validation_name = 'data_batch_2.mat';
8       %X,Y and y data for each file.
9       [X_train,Y_train,y_train] = LoadBatch(training_name);
10      [X_test,Y_test,y_test] = LoadBatch(test_name);
11      [X_validation,Y_validation,y_validation] = LoadBatch(validation_name);
12
13      %Initialize the parameters of the model W and b. [size(W) = 10x3072
14      %size(b) = 10x1]
15      %Initialize each entry to have Gaussian random values with 0 mean and
16      %standard deviation 0.01.
17      c = 0;
18      a = 0.01;
19      K = size(Y_train,1);
20      d = size(X_train,1);
21      b = a.*randn(K,1)+c;
22      W = a.*randn(K,d)+c;
23
24      %n_batch, eta,n_epochs,lambda aprameters
25      lambda = 0;
26      GDparams.eta = 0.1;
27      GDparams.n_batch =100;
28      GDparams.n_epochs =40;
29
30      %Calculate the training and validation loss.
31      Loss_validation = zeros(1,GDparams.n_epochs);
32      Loss_train = zeros(1,GDparams.n_epochs);
33      for j=1: GDparams.n_epochs
34          Loss_validation(j) = ComputeCost(X_validation, Y_validation, W, b, lambda);
35          Loss_train(j) = ComputeCost(X_train, Y_train, W, b, lambda);
36          [Wstar, bstar] = MiniBatchGD(X_train, Y_train, GDparams, W, b, lambda);
37          W=Wstar;
38          b=bstar;
39      end
40      %Plot the training and the validation loss.
41      %Assignment description : around 2.0...
42      figure()
43      plot(1 : GDparams.n_epochs, Loss_train,'g')
44      hold on
45      plot(1 : GDparams.n_epochs, Loss_validation,'r')
46      hold off
47      xlabel('Epochs');
```

```matlab
48      ylabel('Loss value');
49      legend('Training loss', 'Validation loss');
50
51      %After training, just compute the accuracy of your learnt classifer on
52      %the test data. (Assignemnt description = 36.39%)
53
54      P = EvaluateClassifier(X_train,W,b);
55      accuracy_train = ComputeAccuracy(P,y_train);
56      disp(['Training Accuracy:' num2str(accuracy_train) '%'])
57      P = EvaluateClassifier(X_test,W,b);
58      accuracy_test = ComputeAccuracy(P, y_test);
59      disp(['Test Accuracy:' num2str(accuracy_test) '%'])
60
61      %Visualization of the weight matrix as an image.
62      K = 10;
63      for i = 1 : K
64          im = reshape(W(i, :), 32, 32, 3);
65          s_im{i} = (im - min(im(:))) / (max(im(:)) - min(im(:)));
66          s_im{i} = permute(s_im{i}, [2, 1, 3]);
67      end
68      figure()
69      montage(s_im, 'size', [1, K])
70
71      %Do the same with the following configurations:
72      %lambda = 0, n_epochs = 40, n_batch = 100, eta = 0.1
73      %lambda = 0, n_epochs = 40, n_batch = 100, eta =0.01
74      %lambda = 0.1, n_epochs =40, n_batch = 100, eta = 0.01
75      %lambda = 1, n_epochs = 40, n_batch = 100, eta = 0.01
76  end
```