

## COURSEWORK

IMPERIAL COLLEGE LONDON

DEPARTMENT OF BIOENGINEERING

---

# Coursework 2

---

*Author:*

Georges Nomicos (CID: 01571775)

Date: January 5, 2019

1. To choose an appropriate classifier, the size of the available data is considered first. The number of features and distinct class are also taken into account. The neural network method does not require extensive knowledge of the distribution of the data. Neural networks does not make any assumptions about the feature dependence and in contrast captures the relationship between features. Convergence gives an accurate indications when the training is complete. Common issues with neural networks such as overfitting can be solved by using regularization and dropout (or batch normalization) methods, local optima problem can be avoided by using state of the art gradient descent algorithm (Adam, a combination of RMSprop and momentum). Hyperparameters search remains one of the main drawback, where grid search and random search are used (for a chosen range) or Bayesian optimization.

For the training and evaluation of the model the data is shuffled to produce a general model and avoid any biases. The data is then divided in training, validation and testing data (proportion of 80:10:10 respectively). The model parameters will be learned from the training data. While the validation provides in parallel to training the accuracy of the model every chosen iterations. This allows the inclusion of early-stopping which decreases running time and reduce overfitting. The input features are normalized to speed up training time. The mean of each feature vector is subtracted to each sample of the feature vectors and then divided by the standard deviation of each feature vector. It is beneficial later: producing a less elongated gradient descent search and thus the freedom of using higher learning rate. Figure 1 shows the classification pipeline starting from the top with training data which is fed to our neural network and compared to the label at each forward pass. The weights and biases are updated by backpropagating the error and using a gradient descent algorithm to minimize the loss iteratively. Once the model parameters are learned, the latter are stored and can be used to predict the label of new data. The training phase can be computationally expensive but the prediction is done instantly (one forward pass through the network).

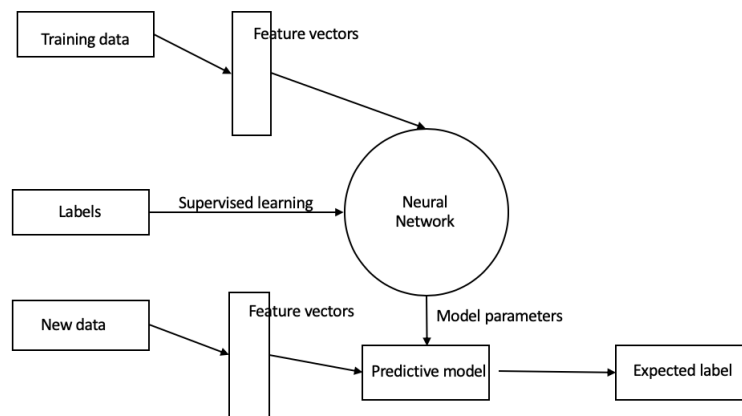


Figure 1: Classification pipeline

2. The classifier chosen is a 5 layer artificial neural network with 80 neurons on the hidden layers. Activation function are ReLU for the first 4 layers and the last one is softmax to produce a readable probabilistic output. Stochastic gradient descent with Adam optimization is implemented to update the weights and biases. Mini-batch is used during training to reduce running time instead of feeding the whole data to the neural network every forward/backward pass. Additionally early-stopping helped reduce overfitting with L2 regularization.
3. The classifier produced an accuracy of 99.5% on the test data set and took 30 seconds to train. As specified early when modelling the neural network the data was divided in three different sets all containing non identical samples, training, validation and testing. Testing set was left to test the performance of our model on unseen data. The validation is used during the training to calculate the accuracy at each forward pass. The latter is used for early stopping and hyperparameters optimization. The training loss was calculated using cross-entropy. The loss was used as an indication of convergence and should approach 0. The stochastic behaviour is caused by the mini-batch algorithm. The validation accuracy was calculated every 20 iterations. Figure 2 shows the validation and training accuracy, and the training loss. Both validation and training curves shows no overfitting and convergence to accuracy above 95% after just 200 iterations (a bit more than 1 epoch). The maximum validation accuracy is reached after 5 epochs.

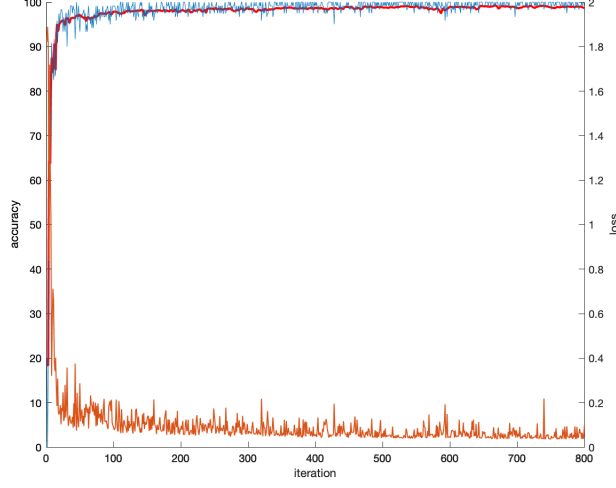


Figure 2: Performance of the classifier (red curve: validation accuracy, blue curve: training accuracy and orange curve: training loss)

A grid search was used to determine the best hyperparameters. The range tested for these hyperparameters are included in Table 1.

Table 1: Hyperparameters optimization range

Hyperparameters	Range
Layers	[2,3,4,5]
Neurons	[30,50,80,100]
Regularization	[ $10^{-5}$ , $10^{-4}$ , $10^{-3}$ , $10^{-2}$ ]
Epoch	[1,3,5,7]
Learning Rate	[ $10^{-4}$ , $10^{-3}$ , $10^{-2}$ , $10^{-1}$ ]

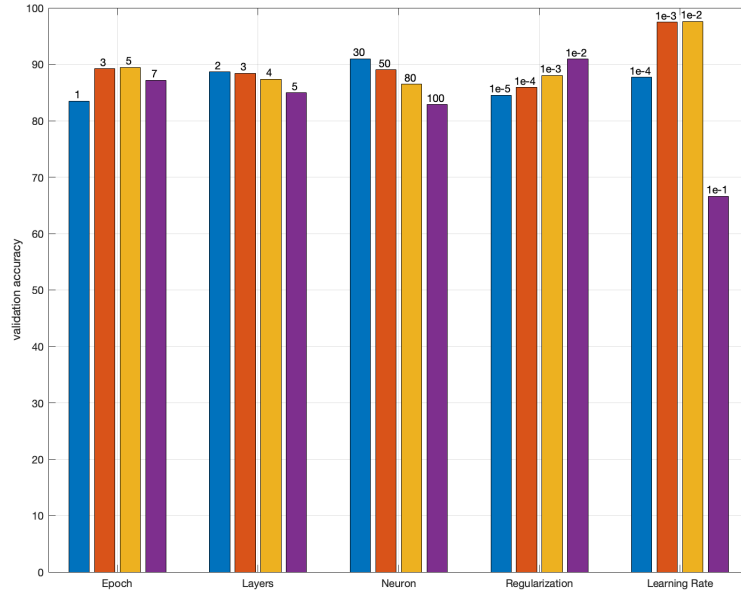


Figure 3: Mean of the validation accuracy of each individual hyperparameter value for all combinations

In total 1024 combinations were tested, taking 2 hours to run. Figure 3 shows the mean of validation accuracy for each individual hyperparameters of each category. The learning rate produces the most

variance in accuracy and mainly affects the accuracy. For the number of epoch, 1 is obviously not enough and 7 is too many (overfit).

The combinations above 99% were selected and are displayed in Table 2. Learning rate and regularization almost all similar. An additional refined grid search could be performed for number of neurons and layers but the accuracy produced is sufficiently high. The last row in Table 3 (Appendix) is chosen for the hyperparameters of our model. Figure 4 contains the accuracies for these hyperparameters.

One of the disadvantages of neural networks is the amount of hyperparameters requiring fine tuning. Out of 8 hyperparameters, 5 were fixed and the others tuned with grid search taking 2 hours. Training time can be time consuming as well, but algorithm like mini-batch and normalizing the data reduces the process. Neural networks cannot be interpreted compared to other methods like decision tree, so it might be difficult or impossible to understand its performance.

When properly implemented and using state of the art algorithm neural network can be very accurate with an accuracy above 99% and training time of 30 seconds. The classification time is considerably shorter time compared to a method like KNN which also requires to store the whole training set for classification. Once hyperparameter are tuned and model is trained, the model predicts the label instantly. Neural networks can fit any data with the drawback of hyperparameters tuning.

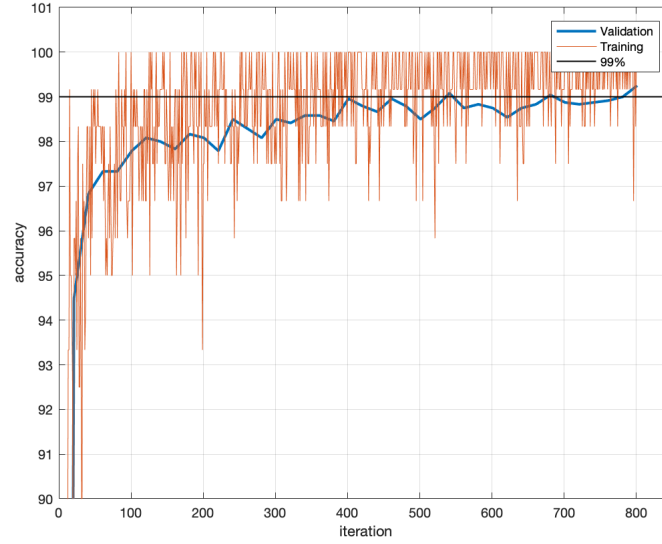


Figure 4: Accuracy of the hyperparameter combination giving the highest validation accuracy

Table 2: Confusion matrix

		Actual class				
		1	2	3	4	5
Predicted class	1	579	0	0	0	0
	2	0	637	2	0	0
	3	0	0	540	1	0
	4	0	0	3	399	9
	5	0	0	0	1	232

## Appendix

Table 3: Combination of hyperparameters resulting a validation accuracy above 99%

Validation Accuracy	Layers	Neurons	Regularization	Epoch	Learning Rate
99.041	5	80	1e-4	7	1e-2
99.041	3	80	1e-4	5	1e-2
99.041	3	80	1e-4	7	1e-2
99.041	4	80	1e-4	7	1e-2
92.082	2	100	1e-4	5	1e-2
99.082	3	50	1e-3	7	1e-2
99.124	2	80	1e-4	7	1e-2
99.166	4	80	1e-3	7	1e-2
99.166	4	100	1e-4	7	1e-2
99.207	5	100	1e-4	7	1e-2
99.249	5	80	1e-4	5	1e-2

## Matlab Code

```
1 %% Train/Vald/Test.
2 %%Training and validation in parallel. Test on the model in the end.
3 clear;
4 clc;
5
6 tic
7
8 load('data.mat');
9 [X,Y,X_V,Y_V,Test_X,Test_Y] = PreProcessing(data,'nequal');
10 X_size = size(X,2);
11 [param,W,B,Ad] = initialization(5,[80 80 80 80 5],X_size,1e-2,...
12 1e-4,0.12,40,5,120); %(Number_of_layer,Neuron_layer,X_size,Learning_rate,...
13 %Regularization,Std_weight,patience,epoch,batchsize) (3,[60 60 5],X_size,1e
14 -3,...
15 %1e-3,0.12,40,38,100)
16 [Loss,accuracy_training,accuracy_val,W,B] = training(X,Y,param,W,B,'adam',Ad,X_V
17 ,Y_V);
18 parameters = struct();
19 parameters.W = W;
20 parameters.B = B;
21 parameters.Number_of_layer = param.Number_of_layer;
22 [error,Conf] = prediction(Test_X,Test_Y,parameters);
23
24
25 %% Hyperparameter optimization
26 clear;
27 clc;
28 tic
29
30 load('data.mat');
31 [X,Y,X_V,Y_V,Test_X,Test_Y] = PreProcessing(data,'nequal');
32 X_size = size(X,2);
33
34 %%Hyperparameters to be chosen:
35 %%Learning rate, regularization factor, std for weight init, number of
36 %%hidden layers, number of neuron per layer, batch size or number of epoch,
37
```

```

38 Neurons = [30 50 80 100];
39 Layers_ = [2,3,4,5]; %for 50 neurons per layer
40 ep = [1,3,5,7];
41 Learning = [1e-4,1e-3,1e-2,1e-1];
42 Regularization_ = [1e-5,1e-4,1e-3,1e-2];
43
44 std_weight = [0.01,0.05,0.1,0.2];
45 batch_size = [50,80,100,150];
46
47 Loss = containers.Map('UniformValues',false);
48 accuracy_training = containers.Map('UniformValues',false);
49 accuracy_val = containers.Map('UniformValues',false);
50
51
52 %Determine how many neuron per layer assuming we have an uniform number of
53 %neuron per layer
54 count = 0;
55 for a = 1:length(ep)
56     epoch = ep(1,a);
57     for b = 1:length(Layers_)
58         Layers = Layers_(1,b);
59         for c = 1:length(Neurons)
60             n_f = Neurons(1,c);
61             N_Layer = {[n_f 5],[n_f n_f 5],[n_f n_f n_f 5],[n_f n_f n_f n_f 5]};
62             %30 50 80 100
63             NN = N_Layer{Layers-1};
64             for d = 1:length(Regularization_)
65                 Regularization = Regularization_(1,d);
66                 for e = 1:length(Learning)
67                     Learning_rate = Learning(1,e);
68                     algo = 'adam';
69                     Std_weight = 0.12;
70                     patience = 30;
71                     batchsize = 120;
72
73                     [param,W,B,Ad] = initialization(Layers,NN,X_size,
74                                             Learning_rate,...
75                                             Regularization,Std_weight,patience,epoch,batchsize); %(
76                                             Number of layer, Number of neuron per layer)
77                     [Loss,accuracy_training,accuracy_val,W,B] = training(X,Y,
78                                             param,W,B,'adam',Ad,X_V,Y_V);
79                     count = count + 1;
80                     disp(count)
81
82                     LossS{count} = Loss;
83                     accuracy_trainingS{count} = accuracy_training;
84                     accuracy_valS{count} = accuracy_val;
85                     HyperPara{count} = {epoch,Layers,n_f,Regularization,
86                                             Learning_rate};
87
88                 end
89             end
90         end
91     end
92 end
93 %% Hyperparameter optimization results processing
94
95 %Search for hyperparameter giving 99% above accuracy
96 count = 0;
97 for k = 1:length(HyperPara)
98     valid = accuracy_valS{k};
99     valid = valid(:,length(valid));

```

```

94     all_v(:,k) = valid;
95     if valid >= 99
96         count = count + 1;
97         index(:,count) = k;
98     end
99 end
100 max_v = max(all_v(:,index)) %Max accuracy
101 index_max = find(all_v == max_v)
102 HyperPara{index_max} %Hyperparameters for max accuracy
103
104 Param = [1,3,5,7; %epoch
105          2,3,4,5; %layers
106          30 50 80 100; %neuron
107          1e-5,1e-4,1e-3,1e-2; %reg
108          1e-4,1e-3,1e-2,1e-1]; %learning
109 sum_v = zeros(5,4);
110
111 Param_s = {'1','3','5','7'; %epoch
112            '2','3','4','5'; %layers
113            '30','50','80','100'; %neuron
114            '1e-5','1e-4','1e-3','1e-2'; %reg
115            '1e-4','1e-3','1e-2','1e-1'};
116
117 for k = 1:5
118     P = Param(k,:);
119     for i = 1:4
120         P_c = P(:,i);
121         count = 0;
122         clear index_all
123         for H_l = 1:length(HyperPara)
124             HP = HyperPara{H_l};
125             HP = cell2mat(HP);
126             index_ = find(HP(:,k)==P_c);
127
128             if isempty(index_)
129                 continue
130             else
131                 count = count + 1;
132                 index_all(:,count) = H_l;
133             end
134         end
135
136         for j = 1:length(index_all)
137             v = accuracy_valS{index_all(:,j)};
138             v_max = v(:,length(v));
139             sum_v(k,i) = sum_v(k,i) + v_max;
140         end
141         sum_v(k,i) = sum_v(k,i)/length(index_all);
142     end
143 end
144
145
146 %Bar chart for individual of each hyperpara in each category
147 hB=bar(sum_v); % use a meaningful variable for a handle array...
148 hAx=gca; % get a variable for the current axes handle
149 hAx.XTickLabel={'Epoch','Layers','Neuron','Regularization','Learning Rate'}; %
150     label the ticks
151 hT=[]; % placeholder for text object handles
152
153 for i=1:length(hB) % iterate over number of bar objects

```

```

154     hT=[hT text(hB(i).XData+hB(i).XOffset-0.01,hB(i).YData,Param_s(:,i), ...
155             'VerticalAlignment','bottom','horizontalalign','center')];
156
157 end
158
159 %Plot max accuracy val adn test
160 figure
161 for l = 1:length(index_max)
162     val = accuracy_valS{index_max(:,l)};
163     tr = accuracy_trainingS{index_max(:,l)};
164     plot(1:20:20*length(val),val,'LineWidth',2)
165     hold on
166     plot(1:length(tr),tr)
167     hold on
168     plot([1 1100],[99 99],'k','LineWidth',1)
169     axis([0 850 90 101])
170 end
171
172 %Display all hyperparameters above 99%
173 for l = 1:length(index)
174     val = accuracy_valS{index(:,l)};
175     val = val(:,length(val))
176     tr = accuracy_trainingS{index(:,l)};
177 end
178 HyperPara{index}
179
180 %Validation and loss for best parameters
181 figure
182 line(1:length(accuracy_val),accuracy_val,'Color','r','LineWidth',2)
183 line(1:length(accuracy_training),accuracy_training)
184
185 ax1 = gca;
186 ax1.XColor = 'k';
187 ax1.YColor = 'k';
188 ax1_pos = ax1.Position;
189 ax2 = axes('Position',ax1_pos,...
190           'YAxisLocation','right',...
191           'Color','none');
192
193 line(1:length(Loss),Loss,'Parent',ax2,'Color','b')
194 %% FUNCTIONS USED ACCROSS
195 function [X,Y,X_V,Y_V,Test_X,Test_Y] = PreProcessing(data,cl_eq)
196 %The string can either be 'equal' or 'nequal' to choose wether or not to have
197 %an equal number of samples for each class. In this function the data is
198 %split in 3 sets (training,validation,testing). Additionnaly the data is
199 %normalized according to the training parameters.
200 label = data(:,1);
201 input = data(:,2:size(data,2));
202
203 %Remove some data to get an equal number of sampling points in each
204 %class, class 5 being the class with the less samples, it is chosen as the
205 %number of samples for each class. It can be commented.
206 if isequal('equal',cl_eq)
207     data = []; % Creates an empty data array
208     class_l = length(label(label==5)); %number of samples for each class
209     start = 1;
210     for c = 1:5 %Shuffle data and picks class_l sample points for each class
211         label_1 = label==c;
212         label_1 = input(label_1,:);
213         remove_1 = size(label_1,1) - class_l;
214         remove = randperm(remove_1);

```



```

215         label_1(remove,:) = [];
216         till = c * class_1;
217         data(start:till,2:65) = label_1;
218         data(start:till,1) = c;
219         start = 1 + till;
220     end
221     label = data(:,1); %label for all classes requires another shuffle
222     input = data(:,2:size(data,2)); %features (64) for each class
223 elseif isequal('nequal',cl_eq)
224 end
225
226 numDatapnts = size(label,1); %Total number of samples used for training/
    validation/testing
227 %The proportion is 80:10:10
228
229 s = RandStream('mt19937ar','Seed',1); %Fix a seed
230 RandStream.setGlobalStream(s)
231 elems = randperm(numDatapnts);
232
233 %To split data in equal sets, we choose 1 and later take a percentage for
234 %the validation and testing
235 n = 1;
236 nDiv = floor(length(elems)/n);
237 start = 1;
238 setsData = zeros(nDiv,n);
239 for j = 1:n
240     till = j*nDiv;
241     till = floor(till);
242     setsData(:,j) = elems(start:till);
243     start = till+1;
244 end
245
246 % 1% of the entire data is selected for training
247 perc = 0.1; %percentage test data
248 test_n = perc * length(elems);
249 test_n = floor(test_n); %Take nearest integer
250 setsData_2 = elems(1:test_n);
251 elems(1:test_n) = [];
252 setsData_1 = elems;
253
254 %Training/Validation data to be split
255 X = input(setsData_1,:);
256 Y = label(setsData_1,:);
257
258 %Calculate mean and standard deviation to normalize the data
259 mean_X = mean(X,1);
260 std_X = std(X);
261 X = X - mean_X;
262 X = X./std_X;
263
264 %Validation data, every time the data is split another shuffling is
265 %applied
266 numVT = length(Y);
267 elems_v = randperm(numVT);
268 perc_v = 0.111; %percentage val data
269 val_n = perc_v * length(elems_v);
270 val_n = floor(val_n);
271 setsData_V = elems_v(1:val_n);
272 elems_v(1:val_n) = [];
273 setsData_T = elems_v;
274

```

```

275 X_V = X(setsData_V,:); %Validation data
276 Y_V = Y(setsData_V,:);
277 X = X(setsData_T,:); %Training data
278 Y = Y(setsData_T,:);
279
280 %Test data which is also normalized by the training mean and std
281 Test_X = input(setsData_2,:);
282 Test_Y = label(setsData_2,:);
283 Test_X = (Test_X - mean_X) ./ std_X;
284 end
285
286 function [param,W,B,Ad] = initialization(Number_of_layer,Neuron_layer,X_size,
    Learning_rate,...
287 Regularization,Std_weight,patience,epoch,batchsize)
288 %Hyperparameters are chosen in this function and the memory is allocated
289 %for the weights,bias and adam parameters matrices. containers.Map is used
290 %to store weights,biases and adam parameters. They can be accessed with a
291 %key (string) and any size cell array can be contained in the same map. It
292 %also makes the code more readable if string are used to access these values.
293 std = Std_weight;
294 param = struct();
295 param.n = Learning_rate; %learning rate
296 param.reg = Regularization; %regularization factor
297 param.epoch = epoch;
298 param.batchsize = batchsize;
299 param.patience = patience;
300 param.Number_of_layer = Number_of_layer;
301 param.Neuron_layer = Neuron_layer;
302 param.beta1 = 0.9; %Fixed hyperp
303 param.beta2 = 0.999; %Fixed hyperp
304
305 Ad = containers.Map('UniformValues',false);
306 Row_w = X_size;
307 W = containers.Map('UniformValues',false);
308 B = containers.Map('UniformValues',false);
309
310 for N = 1:Number_of_layer
311
312     w = strcat('w', num2str(N));
313     b = strcat('b', num2str(N));
314     m = strcat('m',num2str(N));
315     mt = strcat('mt',num2str(N));
316     v = strcat('v',num2str(N));
317     vt = strcat('vt',num2str(N));
318     m_b = strcat('m_b',num2str(N));
319     mt_b = strcat('mt_b',num2str(N));
320     v_b = strcat('v_b',num2str(N));
321     vt_b = strcat('vt_b',num2str(N));
322     Ad(m)=0;Ad(mt)=0;Ad(v)=0;Ad(vt)=0;
323     Ad(m_b)=0;Ad(mt_b)=0;Ad(v_b)=0;Ad(vt_b)=0;
324     W(w) = std * randn(Row_w,Neuron_layer(N));
325     B(b) = zeros(1,Neuron_layer(N));
326
327     Row_w = Neuron_layer(:,N) ;
328 end
329 end
330
331 function [loss,A] = forward_fnc(param,X,Y,W,B)
332 %Forward pass is computed in this function. It takes as input the param
333 %structure, the training input and label, weights and bias containers. It
334 %returns the loss and the activation function container to be fed to the

```

```

335 %backward function.
336
337 %Activation functions for each hidden layer:
338 %ReLU - ReLU - ... - ReLU - Softmax
339
340 N_l = param.Number_of_layer;
341 Neur = param.Neuron_layer;
342 Z = containers.Map('UniformValues',false); %Linear function  $A * X + B$ 
343 A = containers.Map('UniformValues',false); %  $A = f(Z)$  where  $f$  is the activation
    function
344
345     for N = 1:(N_l-1)
346         z = strcat('z', num2str(N)); %z = w*a
347         a = strcat('a', num2str(N)); %a = f(z)
348         w = strcat('w', num2str(N));
349         b = strcat('b', num2str(N));
350
351         Z(z) = X * W(w) + B(b); %  $X = a$  where  $a_0$  is the training set
352         a_to_be = Z(z);
353         a_to_be(a_to_be <= 0) = 0; %ReLU function
354         A(a) = a_to_be;
355
356         X = A(a); %for next iteration
357     end
358
359
360 %Last layer, softmax instead of ReLU
361 z = strcat('z', num2str(N+1)); %z = w*a
362 a = strcat('a', num2str(N+1)); %a = f(z)
363 w = strcat('w', num2str(N+1));
364 b = strcat('b', num2str(N+1));
365
366 Z(z) = X * W(w) + B(b); %N+1 is the last layer
367
368 inter = exp(Z(z));
369 A(a) = inter./sum(inter,2);
370
371 y_hat = A(a);
372 tmp = y_hat(sub2ind([length(Y) Neur(:,N_l)],(1:numel(Y))',Y(:))); %find the
    probability of
373 %the correct class
374
375 W_all = values(W);
376 W_all = cellfun(@(x)x.^2,W_all,'UniformOutput',false); %square all elements of
    each weight matrix
377 W_all = sum(cellfun(@(x) sum(x(:)),W_all)); %sum all elements of each weight
    matrix
378
379 loss = sum(-log(tmp))/length(Y) +param.reg*0.5*W_all;
380 end
381
382
383 function [dW,dB] = backward_fnc(X,Y,A,B,W,param)
384 %Backpropagation function to calculate the gradient.
385
386 N_l = param.Number_of_layer;
387 dW = containers.Map('UniformValues',false);
388 dB = containers.Map('UniformValues',false);
389 dw = strcat('dw', num2str(N_l));
390 db = strcat('db', num2str(N_l));
391 a = strcat('a',num2str(N_l));

```

```

392 a_p = strcat('a', num2str(N_l-1));
393 w_r = strcat('w', num2str(N_l));
394 b_r = strcat('b', num2str(N_l));
395
396 delta_k = A(a); %f(z) for last layer
397 delta_k(sub2ind(size(delta_k), (1:numel(Y)), Y(:))) = delta_k(sub2ind(size(
    delta_k), (1:numel(Y)), Y(:)))-1; %(y_hat-y)
398 delta_k = delta_k/length(Y); %divided by the number of samples
399 dW(dw) = transpose(A(a_p)) * delta_k + param.reg*W(w_r);
400 dB(db) = sum(delta_k,1) + param.reg*B(b_r);
401
402 k = N_l - 1;
403
404 for N = 1:(N_l-2)
405 dw = strcat('dw', num2str(k));
406 db = strcat('db', num2str(k));
407 a = strcat('a', num2str(k));
408 a_p = strcat('a', num2str(k-1));
409 w = strcat('w', num2str(k+1));
410 w_r = strcat('w', num2str(k));
411 b_r = strcat('b', num2str(k));
412
413
414 delta = delta_k * transpose(W(w));
415 delta(A(a)<=0) = 0;
416 dW(dw) = transpose(A(a_p)) * delta + param.reg*W(w_r);
417 dB(db) = sum(delta,1) + param.reg*B(b_r);
418
419 k = k - 1; %to store the gradient decreasingly
420 delta_k = delta;
421 end
422
423 %First layer backprop
424 dw = strcat('dw', num2str(1));
425 db = strcat('db', num2str(1));
426 w = strcat('w', num2str(2));
427 a = strcat('a', num2str(1));
428 w_r = strcat('w', num2str(1));
429 b_r = strcat('b', num2str(1));
430
431 delta = delta_k * transpose(W(w));
432 delta(A(a) <=0) = 0;
433 dW(dw) = transpose(X) * delta + param.reg*W(w_r);
434 dB(db) = sum(delta,1) + param.reg*B(b_r);
435 end
436
437 function [Loss, accuracy_training, accuracy_val, W, B] = training(X, Y, param, W, B,
    update, Ad, X_val, Y_val)
438 %Training function which iterate over a specified number of epoch. Takes as
439 %input the training input and label, param struct, the initial weight and
440 %bias, the gradient descent method to be used (update either 'sgd' or
441 %'adam'), the adam parameters and the validation set. Returns the loss,
442 %accuracy for both training and validation, weights and biases.
443 %Mini-batch are used instead of feeding all the training data at each
444 %forward/backward pass.
445
446 %epoch = batchsize * iteration/12000; %one epoch is one full sweep through all
    the data
447 iteration = (param.epoch*length(Y))/param.batchsize;
448 X_ini = X; Y_ini = Y;
449 Number_of_layer = param.Number_of_layer; %Unroll variables from structure

```

```

450 n_patience = 0;
451 prev = struct(); %stores W B for the previous iteration in case of early
    stopping
452 count = 1;
453
454 for it = 1:iteration
455     % Mini-batch
456     shuffle_indexes = randperm(size(X,1));
457     shuffle_indexes = shuffle_indexes(1:param.batchsize);
458     X_batch = X(shuffle_indexes, :);
459     Y_batch = Y(shuffle_indexes);
460     X(shuffle_indexes,:) = [];
461     Y(shuffle_indexes,:) = [];
462
463     if size(X,1) < param.batchsize %When no batch can be extracted from
464         %the total data set, start a new epoch
465         X = X_ini; Y = Y_ini;
466     end
467
468     [loss,A] = forward_fnc(param,X_batch,Y_batch,W,B);
469     [dW,dB] = backward_fnc(X_batch,Y_batch,A,B,W,param);
470
471     Loss(:,it) = loss;
472     parameters = struct();
473     parameters.Number_of_layer = Number_of_layer;
474     parameters.W = W;
475     parameters.B = B;
476
477     %Run this every 100 iterations
478     if it == 1
479         [error,~] = prediction(X_val,Y_val,parameters);
480         accuracy_val(:,1) = error;
481     elseif floor(it/1) == it/1
482         count = count + 1;
483         [error,~] = prediction(X_val,Y_val,parameters);
484         accuracy_val(:,count) = error;
485
486         % Early Stopping
487         if count == 1 || count == 2 %Does not work for first 2 iterations
488             continue
489         elseif n_patience == param.patience
490             W = prev.W;
491             B = prev.B;
492             break
493         elseif n_patience == 0 && accuracy_val(:,count-1) == error &&
494             accuracy_val(:,count-1) > accuracy_val(:,count-2)
495             n_patience = 1;
496             %Early stopping if the error at this
497             %iteration is lower than previous
498         elseif n_patience > 0 && accuracy_val(:,count-1) == error
499             n_patience = n_patience + 1;
500         else
501             n_patience = 0;
502         end
503     end
504
505     a = strcat('a', num2str(Number_of_layer));
506     [p,y_train] = max(A(a),[],2);
507     error_t = Y_batch - y_train;
508     accuracy_training(:,it) = (1-length(error_t(error_t ~= 0))/length(Y_batch))
        *100;

```

```

508
509     if it == 1 %If the ost is three times the previous then the algorithm is
510         stopped.
511         continue
512     elseif loss >= 8 * Loss(:,it-1)
513         disp('Cost exploded')
514         break
515     end
516
517 %Early-stopping to save time. Patience is defined above. If accuracy
518 %value is continuously repeated n-patience time then the algo is
519 %stopped and returns the previous weight,bias.
520
521 prev.W = W;
522 prev.B = B;
523
524 for ng = 1:Number_of_layer %Update the weights and biases for each layer
525     w = strcat('w', num2str(ng));
526     b = strcat('b', num2str(ng));
527     dw = strcat('d',w);
528     db = strcat('d',b);
529
530     if isequal(update,'adam') %Adam
531         m = strcat('m',num2str(ng));
532         mt = strcat('mt',num2str(ng));
533         v = strcat('v',num2str(ng));
534         vt = strcat('vt',num2str(ng));
535
536         Ad(m) = param.beta1.*Ad(m) + (1-param.beta1).*dW(dw);
537         Ad(mt) = Ad(m) ./ (1-param.beta1.^ng);
538         Ad(v) = param.beta2.*Ad(v) + (1-param.beta2).*(dW(dw).^2);
539         Ad(vt) = Ad(v) / (1-param.beta2.^ng);
540         W(w) = W(w) - param.n * Ad(mt) ./ (sqrt(Ad(vt)) + 1e-8);
541
542         m_b = strcat('m_b',num2str(ng));
543         mt_b = strcat('mt_b',num2str(ng));
544         v_b = strcat('v_b',num2str(ng));
545         vt_b = strcat('vt_b',num2str(ng));
546
547         Ad(m_b) = param.beta1.*Ad(m_b) + (1-param.beta1).*dB(db);
548         Ad(mt_b) = Ad(m_b) ./ (1-param.beta1.^ng);
549         Ad(v_b) = param.beta2.*Ad(v_b) + (1-param.beta2).*(dB(db).^2);
550         Ad(vt_b) = Ad(v_b) / (1-param.beta2.^ng);
551         B(b) = B(b) - param.n * Ad(mt_b) ./ (sqrt(Ad(vt_b)) + 1e-9);
552
553     elseif isequal(update,'sgd') %Stochastic gradient descent
554         W(w) = W(w) - param.n * dW(dw);
555         B(b) = B(b) - param.n * dB(db);
556     end
557
558 end
559
560 if mod(it,100) == 0 %Display loss every 100 iteration
561     disp(loss)
562 end
563 end
564 end
565
566 function [error,Conf] = prediction(input,label,parameters)
567 %function to calculate error between input and label for given weight and

```

```

568 %bias. Returns the error between predicted label and true label as well as
569 %the confusion matrix.
570
571 W = parameters.W;
572 B = parameters.B;
573 N_l = parameters.Number_of_layer;
574 X = input;
575
576 for N = 1:(N_l-1)
577     w = strcat('w', num2str(N));
578     b = strcat('b', num2str(N));
579
580     Z = X * W(w) + B(b); % X = a where a0 is the training set
581     a_to_be = Z;
582     a_to_be(a_to_be <= 0) = 0;
583     A = a_to_be;
584
585     X = A; %for next iteration
586 end
587
588 %Output layer
589 w = strcat('w', num2str(N+1));
590 b = strcat('b', num2str(N+1));
591
592 Z = X * W(w) + B(b); %N+1 is the last layer
593
594 inter = exp(Z);
595 A = inter ./ sum(inter, 2);
596 [p, ypred] = max(A, [], 2);
597
598 error = label - ypred;
599 error = (1 - length(error(error ~= 0)) / length(label)) * 100;
600 % disp(error)
601
602 %confusion matrix %comment this when training
603 Conf = zeros(5, 5);
604 for L = 1:length(ypred)
605     Conf(ypred(L), label(L)) = 1 + Conf(ypred(L), label(L));
606 end
607 % disp(Conf)
608 end

```