

Research Summary July 2020

LN Pandey
Rahul Vashisht

July 20, 2020

1 Introduction

Machine learning is a study of algorithms that learns from structured data. Deep learning is a sub field of machine learning that can learn from unstructured data such as images, videos, and text. These learning algorithms have garnered a lot of interest recently due to large number of applications in multiple domains such as medicine, security, and automation. Deep learning has been hugely successful in the limited tasks such as image recognition, image segmentation, image captioning, and in natural language processing. However, there is lack of proper understanding of reasons behind success or potential failure points.

A subclass of deep models known as attention networks have gained much interest recently in tasks such as image captioning, video summarization, and machine translation. However, there is very little theoretical understanding of these models. Large number of applications of attention networks motivate us to understand these models. An intuitive understanding of attention networks, will enable us to give a better learning algorithm for these models.

2 Attention Model on Mosaic image

Attention networks (similar to attention mechanism) a particular subclass of deep networks have gathered a decent amount of success in various areas such as image captioning, and video captioning. Neural attention is defined as the use of weighted average of instances (low dimensional embeddings), where the neural network determines weights. In particular, a clear understanding of an intermediate layer is possible in these types of networks in a way that is not possible with other standard deep networks. This perspective of the intermediate layer enables some recent tools developed for the apprehension of deep networks and interpolating classifiers, which can learn even in the presence of random noise as well as with systematic corruption. Following sections define a setup and experiments for attention models.

2.1 Experiment setup and details

Let X_i denote the i^{th} class in CIFAR10.

F denote the set of Foreground Classes such that $F = \{X_0, X_1, X_2\}$

B denote the set of Background Classes such that $B = \{X_3, X_4, X_5, X_6, X_7, X_8, X_9\}$

i^{th} Mosaic Image (Data Point) be $M_i = (x_i, y_i)$

where x_i be set of images such that $x_i = \{I_k : 1 \leq k \leq 9\}$ and $y_i = Class(I_1)$

Images I_1 to I_9 are randomly shuffled, sampled such that $Class(I_1) \in F$ and $Class(I_j) \in B$, where $2 \leq j \leq 9$.

2.2 Model

Since the dimension of one Image is $3 \times 32 \times 32$, therefore,

$f(x)$ be the function representing "Focus" network such that

$$f : \mathbb{R}^{3 \times 32 \times 32} \rightarrow \mathbb{R}$$

$g(x)$ be the function representing "Classification" network such that

$$g : \mathbb{R}^{3 \times 32 \times 32} \rightarrow F,$$

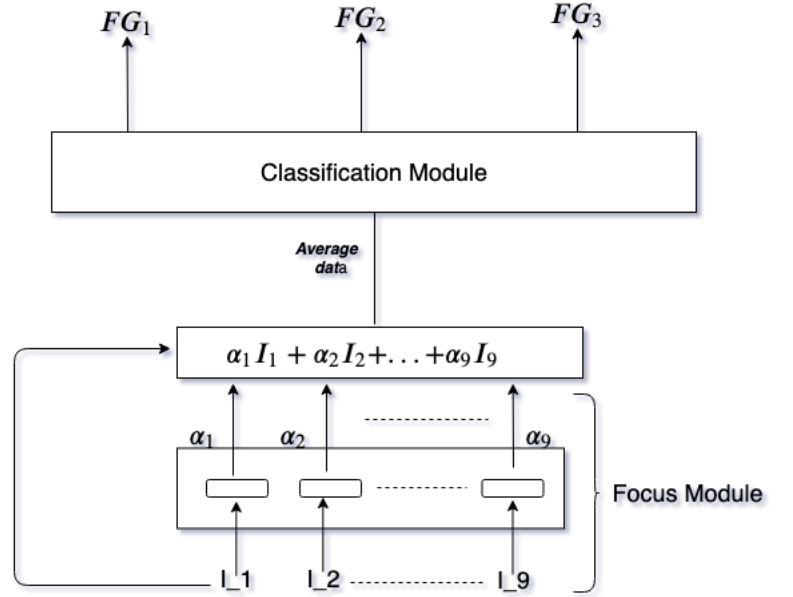
where $F \in [0, 1, \dots, n]$ and n is the number of classes

Let $h(x)$ be the function representing "Model" such that:

$$h : \mathbb{R}^{3 \times 32 \times 32 \times 9} \rightarrow \mathbb{R}^3$$

Let \hat{y}_i be the class predicted by the attention network, then

$$\hat{y}_i = h(x_i) = g\left(\sum_{i=1}^9 I_i f(I_i)\right)$$



3 Experiments done with the above settings:

3.1 Averaging at different Convolution Layers

We perform the experiments by averaging input data at different convolution layers to understand how the performance of attention model varies when we average at different convolution layers. Intuitively, as we average the data at deeper layers the performance should get better. Following table shows the results corresponding to the above experiment.

Averaging at Conv layer	zero layer	first layer	second layer	third layer
training accuracy	96.68	99.55	99.41	99.83
test accuracy	85.58	81.22	88.25	90.01

Table 1: Performance of attention model with averaging different convolution layers

3.1.1 Observations

- There is a gradual increase in performance as we go deeper in layer for data averaging.

3.2 Testing “classification” Network with Different “Focus” Network

1. First, we have created 20k Mosaic images (x_i, y_i) as defined in section 2.1.

Each Dataset D_i will contain 10k weighted aggregated images (avgImg) generated from Mosaic images. Labels of avgImg will be same as y .

avgImg in Dataset D_i is defined as

$$avgImg(D_i) = \frac{i}{9} I_F + \sum_{I_B} \frac{9-i}{9 \times 8} I_B$$

I_F is foreground Image in Mosaic Image x_i and I_B is background Image in Mosaic Image x_i

We consider this weighted aggregated $avgImg(D_i)$ to be the output of "FOCUS" network.

2. Model (Classification network) is trained on Dataset D_i , and Tested on all Dataset D_j , where $j \neq i, i \in [1, 9], j \in [1, 9]$.
3. We want to see, as the Focus network perform better, i.e from Dataset1 to Dataset9 (Focus network learns to focus better by construction of our Datasets), how the Classification network behaves.

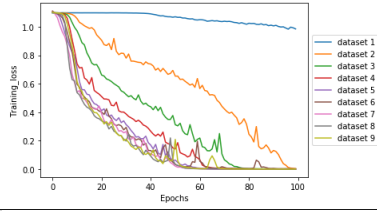
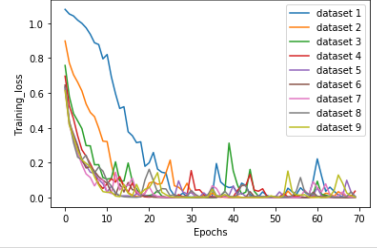
Network used	Train Loss on the Network w.r.t all Datasets
CNN Network	
Mini Inception Network	

Table 2: Train loss on different Network used

3.2.1 OBSERVATIONS

1. Even with $\frac{2}{9}$ focus on foreground images, the classification network is still able to accurately classify 84% of the times.
2. Notice the faster convergence of classification network with more confident focus networks.
3. The classification network is powerful enough to pickup weak signals as well. This can be seen from the significant improvement from dataset2 to dataset3 (and from dataset3 to dataset4).

3.3 Training Classification network by modelling output of Focus network as Dirichlet distribution with different alpha values

3.3.1 Dataset Creation

First, we have created 20k Mosaic images. Labels of every Mosaic image is the class of foreground image present in that image.

Dataset 'i' will contain 10k image. Where each image is weighted average of 9 images present in Mosaic image.

Dataset 1 to Dataset 9 are made from same 10k Mosaic Images.

Dataset 10 is made from the rest 10k Mosaic Images which only contains the true Foreground Image present in every Mosaic Image. We have modeled the output of Focus network using the Dirichlet distribution with parameters alpha Since we know exactly where the foreground image is present in the Mosaic image therefore for Dataset 'i' and $\text{Dir}(\alpha)$,

$$\begin{aligned}
 w_{fg} &= a_i * \alpha \\
 w_{bg} &= \alpha \\
 \text{Avg_Image} &= w_{fg} * I_{fg} + \sum_{bg} w_{bg} * I_{bg}
 \end{aligned}$$

where, a_i is the parameter for Dataset 'i', w_{fg} = weight_of_foreground_image, w_{bg} = weight_of_background_image, I_{fg} = Foreground_Image, I_{bg} = Background_Image

Parameter a_i for different Datasets 'i', can be defined as: $a_i = \left(\frac{i}{9}\right) \left(\frac{72}{9-i}\right)$ for $i \in [1, 8]$.

For Dataset 9, $w_{fg} = \alpha$ and $w_{bg} = 0$.

Dataset 10 is made from the rest 10k Mosaic Images which only contains the true Foreground Image present in every Mosaic Image. We have used $\alpha = [1e^{-2}, 0.5, 2, 1e^5]$
 $\alpha = 1e^{-2}$ denotes sparse samples and $\alpha = 1e^2$ denotes samples with equal weights.

3.3.2 Observation

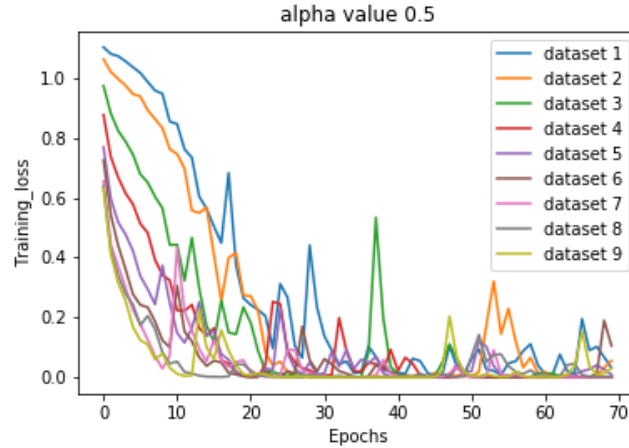


Figure 1: Train Loss for all Datasets when $\alpha = 0.5$

1. With alpha value less than one, occasionally background images will also get focused as input to classification network.
2. These results are similar to the experiments done with randomly fixing the Focus network.

3.4 Stage wise learning

In this experiment instead of training both focus network and classification network simultaneously, we freeze one module, while training other module. Following experiment shows observations when classification network is trained and focus network is freeze at the start for certain number of epochs after which we train other module while freezing the previously trained network.

3.4.1 Observations

Experiment No.	Total_Epochs	Train_network_for_Epoch	What_Learning_Rate	Where_Learning_Rate	Training Accuracy	Testing Accuracy
1	180	5	0.01	0.01	0.96	0.81
2	200	10	0.01	0.01	0.94	0.79
3	280	20	0.01	0.001	0.97	0.43
4	320	20	0.01	0.01	0.97	0.79
5	360	30	0.01	0.01	0.99	0.76
6	500	50	0.01	0.01	0.86	0.41
7	600	50	0.001	0.01	0.97	0.82
8	800	100	0.001	0.01	0.97	0.83

Table 3: Effect of Learning Rate on stage wise training

1. The behaviour of accuracy and CE Loss for alternate training was a bit absurd when train alternatively for relatively larger number of epochs.
2. It seems that a little improvement in Classification network rewards more as compared to a little improvement in Focus network.
3. That is improvement in Classification will lead to more gains!
4. “Classification” Network tries over fits the data while “focus” network was frozen, But when “focus” net starts training, it tries to restrict “classification” net in doing so and tries to learn focus correctly. hence every time after the “Classification” network finishes its training, There is a drop in accuracy and a sudden jump in CE Loss.

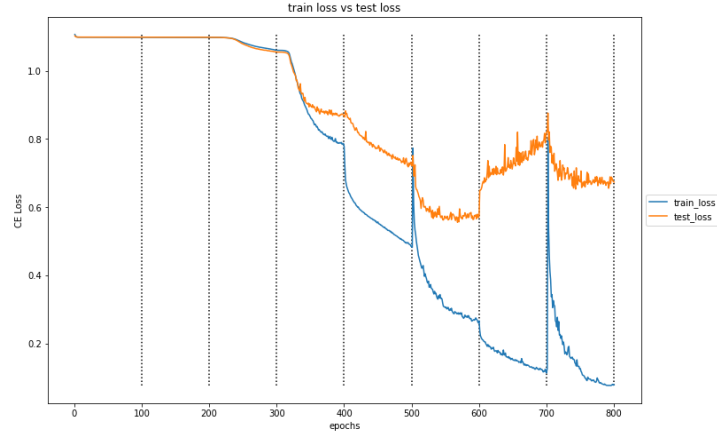


Figure 2: Train Loss and Test loss for every 100 epochs

3.5 Deep Dive into Attention

We have 4 Datasets of Mosaic images (created from CIFAR10 having classes $C = \{X_i : i \in [0, 9]\}$) s.t in Dataset $D_i, F = \{i, i+1, i+2\}, B = C - F$. For all 4 datasets, we tried to analyse the FTPT (Focus True Prediction True), FFPT (Focus False Prediction True), FTPF (Focus True Prediction False), FFPT (Focus False Prediction False) for the following:

1. Data points (M_i s.t $(f, b_0) \in M_i$) having Foreground image f , s.t $\forall f \in F$ and Background image b_0 , s.t $\forall b_0 \in B$.
2. Data points (M_i s.t $(f, b_0, b_1) \in M_i$) having Foreground image f , s.t $\forall f \in F$ and Background images (b_0, b_1) s.t $\forall b_0 \in B$ and $\forall b_1 \in B$.
3. Data points (M_i s.t $(f, b_0, b_1, b_2) \in M_i$) having Foreground image f , s.t $\forall f \in F$ and Background images (b_0, b_1) s.t $\forall b_0 \in B$, $\forall b_1 \in B$ and $\forall b_2 \in B$.
4. Data points (M_i s.t $f \in M_i$ and $(b_0) \notin M_i$) having Foreground image f , s.t $\forall f \in F$ and Background image b_0 , s.t $\forall b_0 \in B$.
5. Data points (M_i s.t $f \in M_i$ and $(b_0, b_1) \notin M_i$) having Foreground image f , s.t $\forall f \in F$ and Background images (b_0, b_1) s.t $\forall b_0 \in B$ and $\forall b_1 \in B$.
6. Data points (M_i s.t $f \in M_i$ and $(b_0, b_1, b_2) \notin M_i$) having Foreground image f , s.t $\forall f \in F$ and Background images (b_0, b_1) s.t $\forall b_0 \in B$, $\forall b_1 \in B$ and $\forall b_2 \in B$.

3.5.1 Observation :

1. Training accuracy may be different for different Foreground classes but For a specific class in foreground Classes, Training accuracy in every case is almost same.
2. Testing Accuracy for every foreground class is almost same.

4 Deep Network Training under random label noise:

4.1 Out of 50k data points 0,5,10,20,30,40,50,100% data points are corrupted

1. On CIFAR10 dataset where the size of dataset is 50k, randomly assign labels to 0,5,10,20,30,40,50,100 percentage of the data and then train the network. Analyse the result for Training and Testing accuracy, CE Loss in all the scenarios.
2. Few observations are:

Corruption Percentage	Train accuracy on corrupted	Train accuracy on un-corrupted	Train accuracy on full	Test accuracy on corrupted	Test accuracy on un-corrupted	Test accuracy on full	Test accuracy on true data
0	-NIL-	1	1	-NIL-	-NIL-	-NIL-	0.83
5	1	1	1	0.0858	0.7830	0.7481	0.78
10	1	1	1	0.1028	0.7509	0.6841	0.75
20	1	1	1	0.1028	0.6659	0.5505	0.66
30	1	1	1	0.1019	0.5864	0.4405	0.58
40	1	1	1	0.1031	0.5058	0.3462	0.50
50	1	1	1	0.0970	0.4252	0.2625	0.42
100	0.99	-NIL-	0.99	0.1084	-NIL-	0.1084	0.10

Table 4: Analysis of Training and Testing Accuracy

4.2 Out of 50k data points, 25k data points were fixed as true data

1. Keep the true training data constant as 25k, add randomise data (corrupted labels or corrupted training data) so that the final corruption percentages are given as 0,5,10,20,30,40,50. Analyse the result for Training and Testing accuracy, CE Loss in all the scenarios.
2. Few observations are:

Corruption Percentage	Train accuracy on corrupted	Train accuracy on un-corrupted	Train accuracy on full	Test accuracy on corrupted	Test accuracy on un-corrupted	Test accuracy on full	Test accuracy on true data
0	-NIL-	1	1	-NIL-	-NIL-	-NIL-	0.83
5	1	1	1	0.0978	0.7496	0.717	0.75
10	1	1	1	0.0844	0.7161	0.651	0.71
20	1	1	1	0.0961	0.6486	0.5353	0.64
30	1	1	1	0.1019	0.5755	0.4325	0.57
40	1	1	1	0.1047	0.5077	0.3479	0.50
50	1	1	1	0.1035	0.4254	0.2658	0.42
100	0.99	-NIL-	0.99	0.1084	-NIL-	0.1084	0.10

Table 5: Analysis of Training and Testing Accuracy

5 Systematic Corruption:

Here we introduce systematic corruption in data by changing a percentage of labels in the data, to understand if neural networks can learn in the presence of high noise percentage. The noise here is not random but we change the image labels systematically as explained in the following section.

5.1 CIFAR10 Dataset

Take CIFAR10 and create training data x_i , z_i as follows. Let x_i be the CIFAR10 images, and y_i be the CIFAR10 labels. The training labels z_i are simply given as follows:

1. $z_i = y_i$ if y_i is either 0 or 1 or 2, z_i is randomly 0 or 1 or 2 otherwise.
2. Train a network on this x_i , z_i data and predict labels on all test images. Check the accuracy on just the 0,1,2 labelled test images.
3. Training was done on CNN (Simple) and Mini-Inception Network (Complex models).
4. Few Observations are:

true training data	Corrupted Training data	Total Training data	Training accuracy	Test accuracy	Test accuracy 0-1-2
1000	35000	36000	1	0.11	0.38
2000	35000	37000	1	0.14	0.49
4000	35000	39000	1	0.16	0.56
6000	35000	41000	1	0.19	0.63
8000	35000	43000	1	0.20	0.67
10000	35000	45000	1	0.21	0.71
12000	35000	47000	1	0.22	0.73
15000	35000	50000	1	0.23	0.78

Table 6: Analysis of CIFAR Corrupted data on Mini-Inception Network

5.2 MNIST Dataset

Take MNIST and create training data x_i , z_i as follows. Let x_i be the MNIST images, and y_i be the MNIST labels. The training labels z_i are simply given as follows:

1. $z_i = y_i$ if y_i is either 0 or 1, z_i is randomly 0 or 1 otherwise.
2. Train a network on this x_i, z_i data and predict labels on all test images. Check the accuracy on 0 and 1 labelled test images.
3. Training was done on CNN (Simple) and Mini-Inception Network (Complex models).
4. Few Observations are:

true training data	Corrupted Training data	Total Training data	Training accuracy	Test accuracy	Test accuracy 0-1
100	47335	47435	1	0.13	0.61
500	47335	47835	1	0.16	0.80
1000	47335	48335	1	0.17	0.83
2000	47335	49335	1	0.19	0.92
4000	47335	51335	1	0.20	0.95
6000	47335	53335	1	0.20	0.96
8000	47335	55335	1	0.20	0.96
12665	47335	60000	1	0.20	0.98

Table 7: Analysis of MNIST Corrupted data on CNN Network

6 Visualising Mosaic Data by Generating 2-Dim Separable Data

Here, we work with 2-dimensional data (to create mosaic data), which is easy to visualise to understand the attention networks. Also following experiments can be used to get an understanding of when the attention networks are interpretable.

6.1 Generation of 10 Classes Using 2 Dimensional data

We sampled 5000 (500 for every class approx) 2-dim data points from multivariate Normal with 10 different Mean vectors and same covariance Matrix. Covariance Matrix being the diagonal matrix with entries [0.1, 0.1].

6.2 Generation of Mosaic Data

1. Available Classes = Class 0, Class 1, Class 2, Class 3, Class 4, Class 5, Class 6, Class 7, Class 8, Class 9.
2. foreground_classes = Class 0, Class 1, Class 2
3. background_classes = Class3, Class 4, Class 5, Class 6, Class 7, Class 8, Class 9
4. Every class will have a 2-dim Data Point. 1 data point was chosen at random from any foreground class, and rest 8 data points are chosen from background classes.
5. Now we have 9 data points which can be arranged randomly in 3×3 image matrix. In particular Dim of Matrix will be 3×6 .

6.3 Observation

Following is the table of parameters and accuracy for different data sets used. Main goal here is to find when the attention networks are interpretable from this experiment.

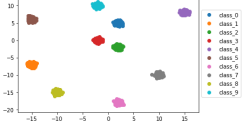
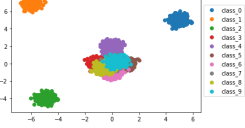
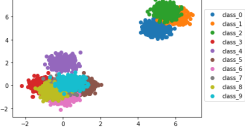
type of dataset used in this	epochs	Focus LR	Classification LR	Train Acc	Test Accuracy
	80	0.01	0.01	100	100
	10	0.01	0.01	100	100
	120	0.01	0.01	88	89

Table 8: Train accuracy and test accuracy

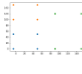
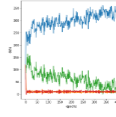
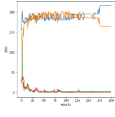
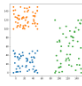
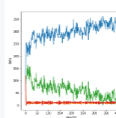
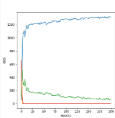
Data Distribution	Linear SVM Accuracy on Mosaic Data	Deep Network Accuracy on Mosaic Data	Linear Attention network Accuracy	Deep Attention Newtork Accuracy	Analysis Linear Attention Network	Analysis Deep Attention Network
	68.22 (not covering)	86.31	90	99.1		
	74.06	91.9	94.6	100		

Figure 3

Now we have repeated the experiments mentioned in table 9 with two dimensional data having two foreground classes and one background class to create mosaic data. The results are shown in the following figure 3

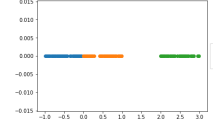
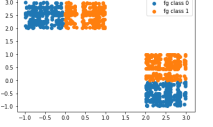
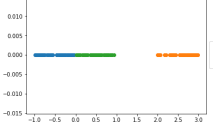
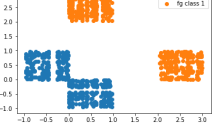
type of data	plot of dataset used	Mosaic Data generated
linearly not seperable		
linearly seperable		

Table 9: Type of Data used for the experiment

1. When mosaic data is not linearly separable The models needs attention to learn a good classifier. For higher dimension once it learns some good attention weights, the attention weight learning stops.
2. When mosaic data is linearly separable The model does not need attention to learn a good classifier. For higher dimension model need not require to learn the attention weights, but is sometimes learning good attention weights.
3. Even with 2 dimensional input data and 6-dimensional mosaic data (not linearly separable), deep attention network performs better than linear attention networks as well as when attention is not used.

7 Correlation of Loss with Foreground and Background Classes

7.1 Generation of 10 Classes Using 2 Dim data

1. We sampled 1000 (100 for every class approx) 2-dim data points from multivariate Normal with 10 different Mean vectors and same Covariance Matrix.
2. Covariance Matrix being the diagonal matrix with entries [0.01, 0.01].

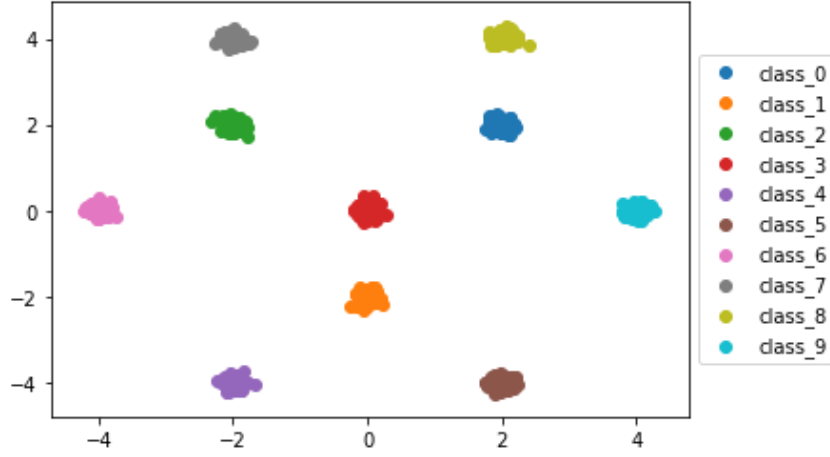


Figure 4: Generation of 10 Classes Using 2 Dim data

7.2 Generation of Mosaic Data

1. Available Classes = Class 0, Class 1, Class 2, Class 3, Class 4, Class 5, Class 6, Class 7, Class 8, Class 9.
2. foreground_classes = Class 0, Class 1, Class 2
3. background_classes = Class3, Class 4, Class 5, Class 6, Class 7, Class 8, Class 9
4. Every class will have a 2-dim Data Point. 1 data point was chosen at random from any foreground class, and rest 8 data points are chosen from background classes.
5. Now we have 9 data points which can be arranged randomly in 3 x 3 image matrix. In particular Dim of Matrix will be 3×6 .

7.3 Model used

Model is developed as combination of 2 modules. "FOCUS" Module learns to "where" the foreground image is present out of 9 images in Mosaic image. "CLASSIFICATION" Module learns "what" is the class of this foreground image out of those 3 foreground classes.

7.4 Input to Model

Mosaic image is input to Module 1 i.e "FOCUS Network", and tries to focus on foreground image present in Mosaic Image. In Particular, Each image (2×1) is input to "FOCUS Network" and hence a 18×1 tensor (9 images) is input to "CLASSIFICATION Network". "FOCUS Network" tries to Focus on Foreground image and returns weighted average of all 9 images. This image is now input to "CLASSIFICATION Network" which finally predicts the Class label of foreground Image.

7.5 Gradient Definition

$$\nabla \alpha_{fg} = \frac{1}{size(Datapoints)} \left[\sum_{datapoints} \nabla Loss_{\alpha}(Fg_one_hot) \right]$$

$$\nabla \alpha_{bg} = \frac{1}{size(Datapoints)} \left[\sum_{datapoints} \nabla Loss_{\alpha}((1 - Fg_one_hot)/8) \right]$$

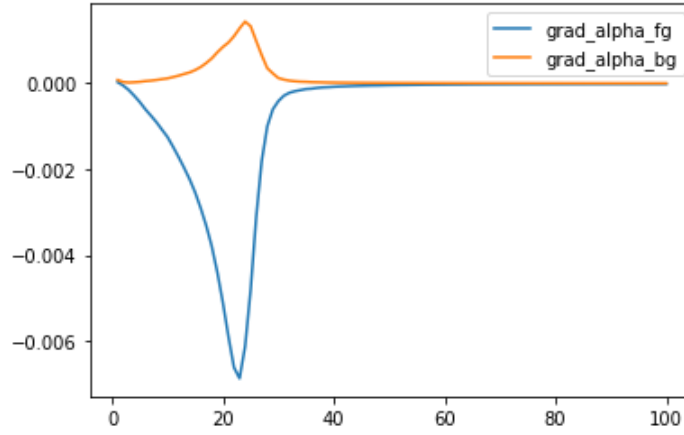


Figure 5: Gradient(Loss) w.r.t alphas for Foreground and Background Classes

8 Confounded Noise

In the following experiments, we introduce label noise to data called as confounding noise to visualise if neural network can learn in the presence of the confounding noise even when data available without noise is very small.

8.1 training Data

Take 2 Dim data and create training data x_i, y_i as shown in table below Training data after Corruption : $z_i = y_i$ if y_i is either 0 or 1. z_i is randomly 0 or 1 otherwise. Train a network on this x_i, z_i data and predict labels on all test images. Check the accuracy on just the zero and one labelled test images.

Data Set	True training data	Corrupted training data
linearly separable at corner		
linearly separable at center		
moon data linearly not separable at corner		

Table 10: True Training data and after corruption

8.2 Observations

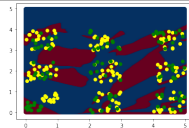
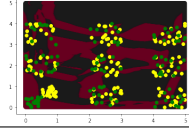
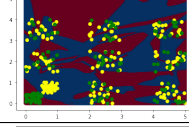
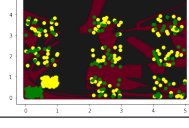
true data	Corrupted Data	Training Accuracy	Test Accuracy	Data Accuracy with no noise	Plots after Training
30	240	97	57	100	
60	240	96	61	100	
90	240	96	66	100	
120	240	98	67	100	

Table 11: training on noisy Linearly separable Data at corner with increasing nos. true data points

Similarly, the same experiments were performed on the moon data linearly separable at corner and data linearly separable at center and following interesting plots after training were observed.

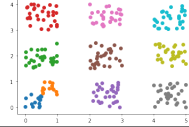
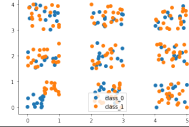
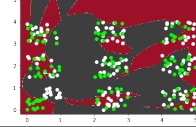
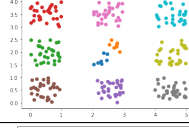
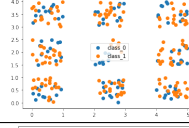
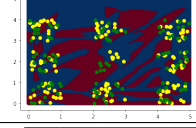
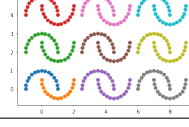
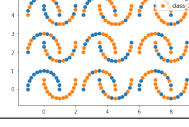
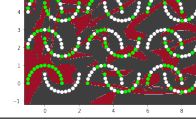
True training data	Corrupted training data	Plots after Training
		
		
		

Table 12: training on Different noisy dataset

Neural Networks can be generalised even in the presence of confounded noise (adding noise to the label). The experiment is similar to the experiment done in the paper "Model Neural Networks Generalize on Small Data sets". Only with small nos. of true data points, model is able to learn the classifier.

9 Averaging Noise

- 3 dim Data points is generated s.t class 0, 1, 2 are in one plane (P_fg), and Class 3, 4, 5, 6, 7, 8, 9 are in other plane (P_bg). P_bg is fixed as x-y plane. Class 0, 1, 2 are Foreground classes, Class 3, 4, 5, 6, 7, 8, 9 are background classes.

2. Generation of Datasets

- First, we have created 1000 Mosaic images (x_i, y_i) as defined in ??.

Each Dataset D_i will contain 1000 weighted aggregated images (avgImg) generated from Mosaic images. Labels of avgImg will be same as y .

avgImg in Dataset D_i is defined as

$$avgImg(D_i) = \frac{i}{9} I_F + \frac{9-i}{9 \times 8} I_B$$

I_F is foreground Image in Mosaic Image x_i and I_B is background Image in Mosaic Image x_i

(b) Model (Classification network) is trained on Dataset D_i , and Tested on all Dataset D_j , where $j \neq i, j \in [1, 9]$.

(c) We want to see, as the Focus network perform better, i.e from Dataset1 to Dataset9 (Focus network learns to focus better), how the Classification network behaves.

3. following are the observation with different angle between the planes. A small plane denote the plane of FG classes.

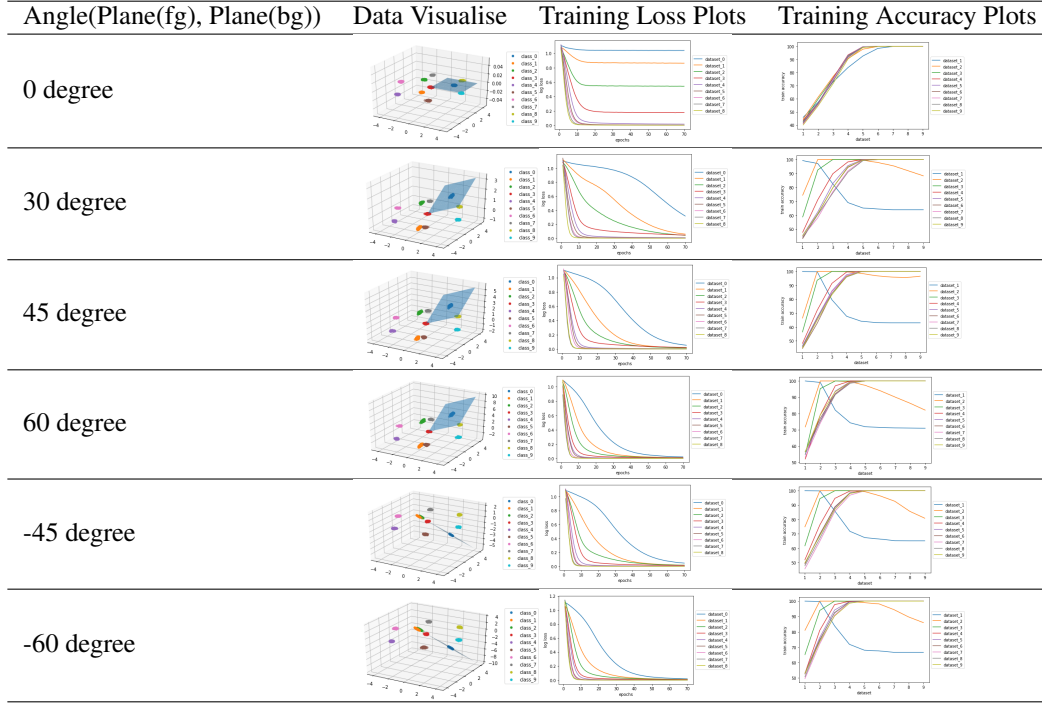


Table 13: training on Different noisy dataset

9.0.1 Observations :

1. As the angle between Foreground classes plan and Background classes plane increases, The loss plots are more steeper.
2. For the trained model, the capability of differentiating between foreground classes and background classes increases with an increase in the angle between the planes.