

CSCI 5922 -Neural Networks and Deep Learning

Lab - 2 Solutions

Gowri Shankar Raju Kurapati
Student ID 110568555

September 26, 2022

1 Influence of Regularization

1.1 Dataset & Hyperparameters Used

- CIFAR 10 Dataset
 - Source : From *torchvision.datasets*
 - The CIFAR-10 dataset used consists of 50000 32x32 color images in 10 classes, with 5000 images per class.
 - Since these are color images, each image has the tensor dimension of $3 * 32 * 32$ where the first dimension is RGB colors.
 - The 10 classes (labels) are *airplane*(0), *automobile*(1), *bird*(2), *cat*(3), *deer*(4), *dog*(5), *frog*(6), *horse*(7), *ship*(8), *truck*(9)
- Hyperparameters
 - training/validation/test split is 70 : 15 : 15
 - The input size to the neural network is $3 * 32 * 32$ image.
 - As stated above, the output layer consists of 10 neurons to predict 10 classes of the CIFAR 10 dataset.
 - The learning rate is set to 0.001 with a batch size of 64 and is held constant for the experiment.
 - The models are trained for 10 epochs.
- Methods:
 - All (Eight) neural networks are Convolution Neural Networks with a series of convolution blocks. Each convolution block consists of a
 - * conv2d layer with max pool
 - * followed by a conv2d layer, max pool and batchNorm (if it is opted as regularizer)
 - Two CNN architectures are used
 - * Two Convolution blocks followed by three fully connected Layers

- * Three Convolution blocks followed by three fully connected Layers.
- * Convolution Block has been described in the previous bullet point.
- I have considered three regularizers
 - * Batch Normalization
 - * Drop out at Fully Connected Layers
 - * Training with Data Augmentation.
- Loss function used for the experiment is the Cross Entropy.
- Adam Optimizer with a learning rate of 0.001 is used.
- I have iterated the experiment with either of two regularizers (Batch Norm & Drop Out) and no regularizer with two architectures defined above and on the training dataset specified (Without any augmentation) resulting in six models.
 - * The models are named Lab2_P1_Regularizer_<RegularizerUsed/NoRegularizer>_<2/3>LayerCNN to capture the number of the convolution blocks and regularizer used.
- I have also iterated with no regularizer but training the models by augmenting the training dataset during training on two architectures resulting in two different models.
 - * The models are named as Lab2_P1_Regularizer_DataAugmentation_<2/3>LayerCNN.
 - * The Augmentation used is a Random Horizontal Flip followed by random cropping of the image with different paddings.
- All the models are trained on Google Collab using the single GPU instance.

1.2 Results & Analysis

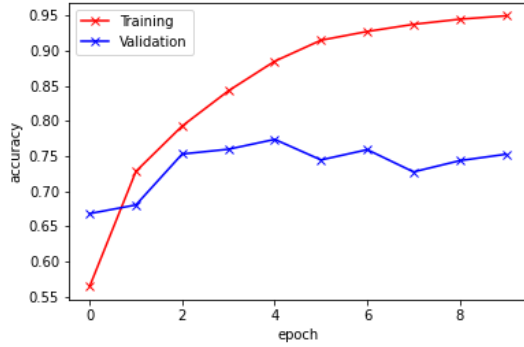
Reporting **Training Times** in seconds for the Eight CNN Models for comparison

CNN Block Layers / Regularizer	Batch Norm	Drop Out	Data Augmen- tation	No Regularizer
2	106.56	103	136.60	109.15
3	115.19	108.6	147.88	118.6

- From the table above, we could infer that,
 - [**3 Layer Networks and DataAugmentation Take more Time to Train**] When the 2 Layer CNN training times are compared with three layers CNN of the same regularizer, the time has increased significantly with an average difference of 9 seconds. Since the number of epochs is constant for all the models, the 3-layer network has to update significantly more model parameters compared to the corresponding 2 Layer. Also, among all the different variations of regularizer used, training with data augmentation takes significantly more time. This can be because of the augmentation taking place during the batch training. The trainLoader being used to generate the batch of images for the training is augmented at the runtime of the train loader for the batch. This processing time is included during the training and hence a significant increase.

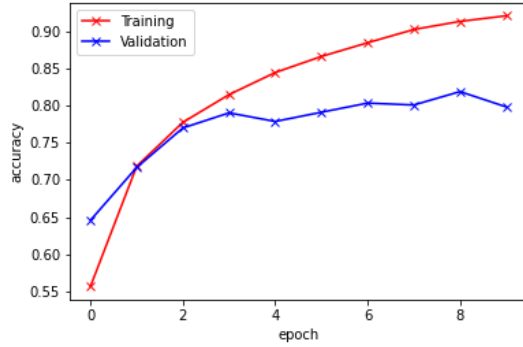
Reporting **Final Validation Accuracies** and **Learning Curves** for the Eight CNN Models for comparison

Acc vs epochs for model Lab2_P1_Regularizer_BatchNorm_2LayerCNN



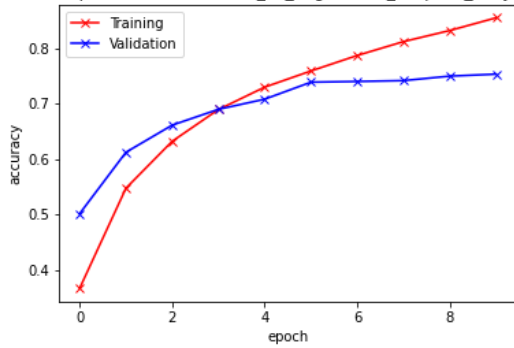
(a) Batch Norm - 2 Layer CNN

Acc vs epochs for model Lab2_P1_Regularizer_BatchNorm_3LayerCNN



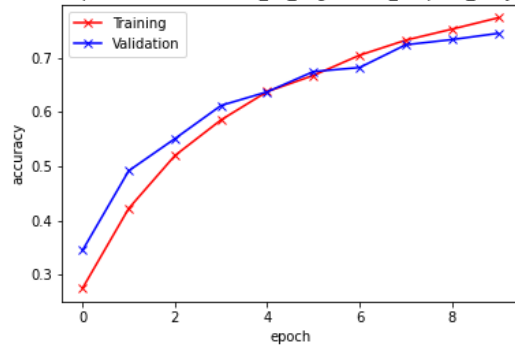
(b) Batch Norm - 3 Layer CNN

Acc vs epochs for model Lab2_P1_Regularizer_DropOut_2LayerCNN



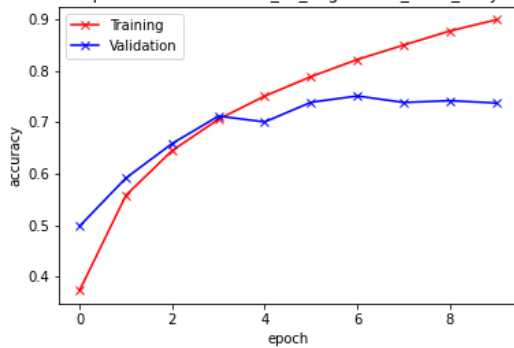
(a) Drop Out - 2 Layer CNN

Acc vs epochs for model Lab2_P1_Regularizer_DropOut_3LayerCNN



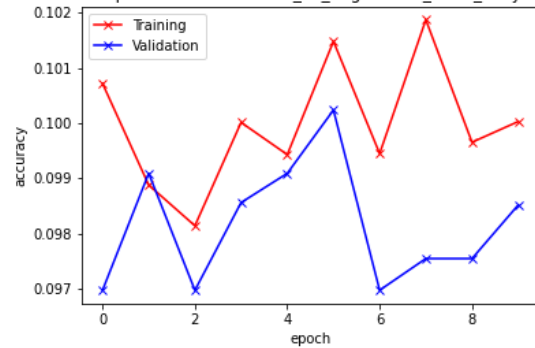
(b) Drop Out - 3 Layer CNN

Acc vs epochs for model Lab2_P1_Regularizer_None_2LayerCNN



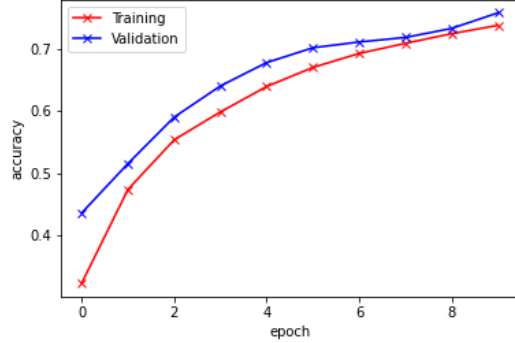
(a) No Regularizer - 2 Layer CNN

Acc vs epochs for model Lab2_P1_Regularizer_None_3LayerCNN

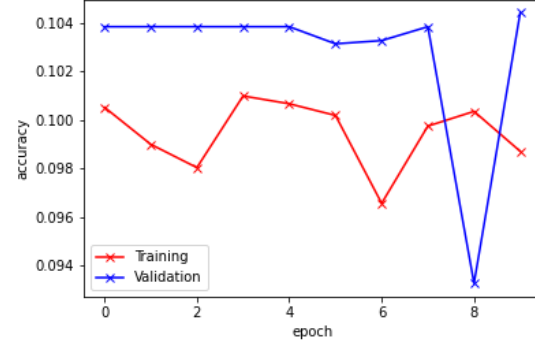


(b) No Regularizer - 3 Layer CNN

Accuracy vs epochs for model Lab2_P1_Regularizer_DataAugmentation_2Layer vs epochs for model Lab2_P1_Regularizer_DataAugmentation_3Layer



(a) Data Augmentation - 2 Layer CNN



(b) Data Augmentation - 3 Layer CNN

CNN Block Layers / Regularizer	Batch Norm	Drop Out	Data Augmen- tation	No Regularizer
2	73.97	75.49	76.23	74.08
3	79.07	75.08	9.94	9.76

- From the above table we can infer that
 - [2 Layer CNN - Data Augmentation works Better]** When no regularizer is used, we see from the corresponding learning curve that there is a little overfitting as the validation accuracy gets stagnant but training accuracy increased. Though using batch normalization doesn't significantly increase the performance, it overfits more. It may be because the model is not able to learn any new features. This assumption is true because once we use Data Augmentation to train thereby making the model more robust to different angles of the image, we see that overfitting is gone (from the corresponding learning curve) and accuracy is increased by nearly 2.5 %.
 - [3 Layer CNN]** We clearly observe that the from the learning curves of Data Augmentation 3 Layer and No Regularizer, the accuracies are very low and there is no pattern to observe. This is expected as we are trying to learn a lot of features (model parameters compared to 2Layer CNN) in just 10 epochs. Since no regularizer is used and on top of it we are also augmenting data with an intention of learning robust features but not giving enough epochs for the model to train as converging takes more time, and a lot of epochs are needed to make proper inferences. Batch Normalization gives more accuracy as normalizing the features at every convolution layer block will ease the updates to move faster in the loss space, thereby decreasing the number of epochs for convergence.
 - [Comparing 2 Layer - 3 Layer CNN for a regularizer]** Since we have established that number of epochs wasn't sufficient for 3 Layer Data Augmentation and NO regularizer, I am omitting the comparison for the corresponding 2 Layer network. We can then observe that the overfitting decreased when the number of layers is increased (so are the accuracies) for the network using the same regularizer. This can be explained as the model capacity has been increased to accommodate more features to learn.
- Thus, the top performing model is 3 Layer CNN with Batch Normalization (Lab2_P1_Regularizer_BatchNorm_3LayerCNN). The Test Accuracy achieved when it is trained on both training split and validation split is **80.81%** .

2 Interpreting CNN Representations

2.1 Network Used And Interpretations used

3 Layer CNN with Batch Normalization (Lab2_P1_Regularizer_BatchNorm __3LayerCNN)

- We use the best model from question-1 to interpret its representations.

[**Interpretation Used**] Retrieving images with similar feature representations. The Feature representations are taken on two spots.

- At the end of the second Convolution Layer Block, when flattened gives you 8192 length vector.
- At the output Layer (at the end of the last fully connected Layer), when flattened has 10 length feature vector

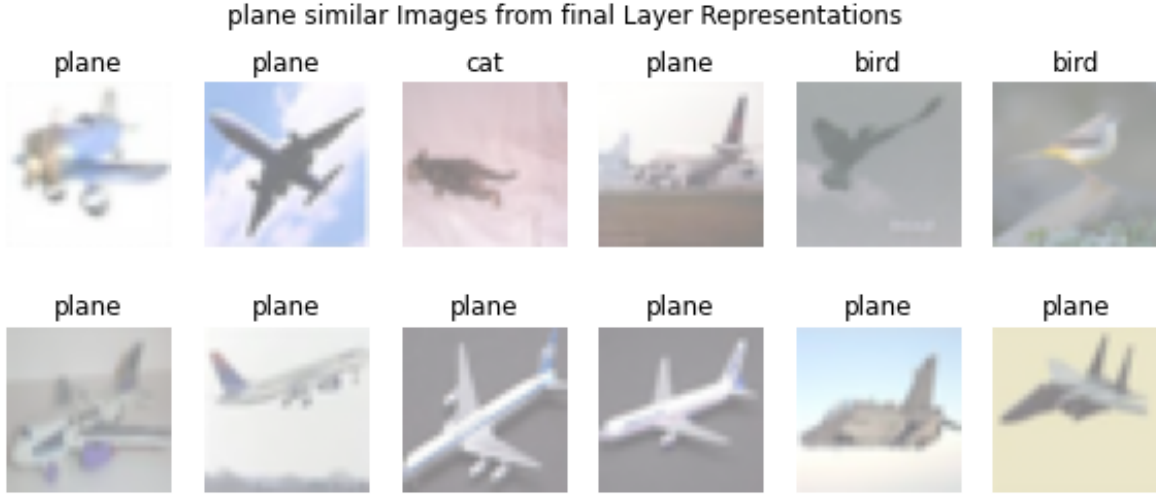
[**Method Used**] For an image in each class, I have taken the representation vector and did a cosine similarity with all the images of the test Dataset and displayed the top 10 results with the highest similarity. The intuition behind this is that at a certain level, images having certain features will be nearer to each other in a feature space. By looking at the images, we can understand what features are more prominent at the end of each layer.

2.2 Results & Analysis

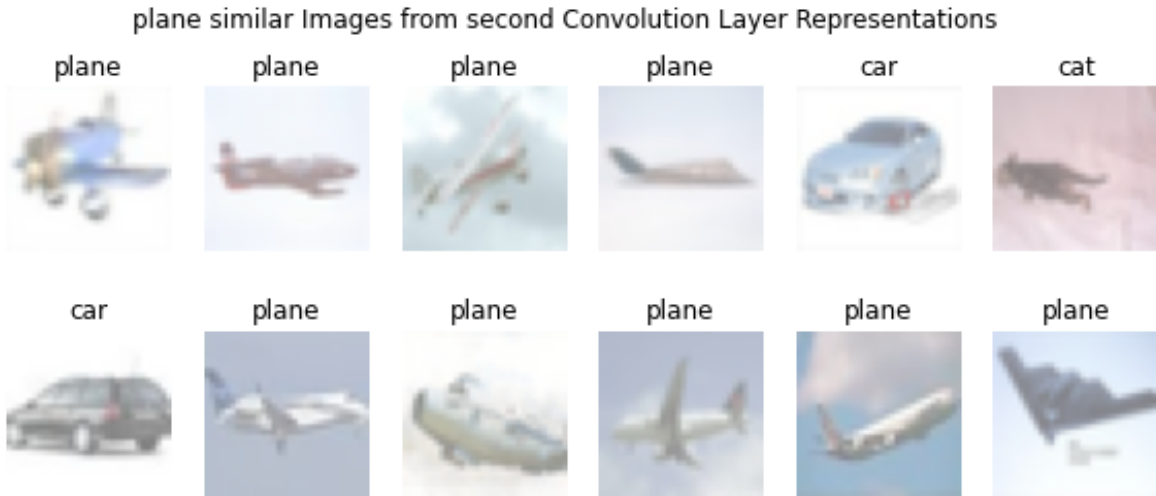
- [**PLANE SIMILAR IMAGES**] For instance, if we retrieve similar images of a plane (*figure-1 (a) & (b)*), we can see that when convolution layer 2 feature representations are used (*figure-1(b)*) the similar images are mostly planes stating that the classification is being done with good accuracy by this layer. But we can also see observe that there are images of cars and cat. Since these are the top 10 similar images, the high-level feature might be an image with a sharp nose, as we can see the car image is like protruding nose tip, and also for the cars, the front portion is like a plane's tip/nose.

Also, when the same level comparison is made with representations of the end layer, we can see all the images are plane except the same cat image we talked about earlier. Here, we don't see any cars as the model would have learned what a plane should be like at this point (As this is the last layer) but it still misclassifies the cat and a bird as planes. The model should be seeing a tip/nose, wings, and a tail. As both the cat image and bird image have all three, the model might be misclassifying it.

Figure 1: Similar Images for the "first Image of first row" [PLANE] in each section



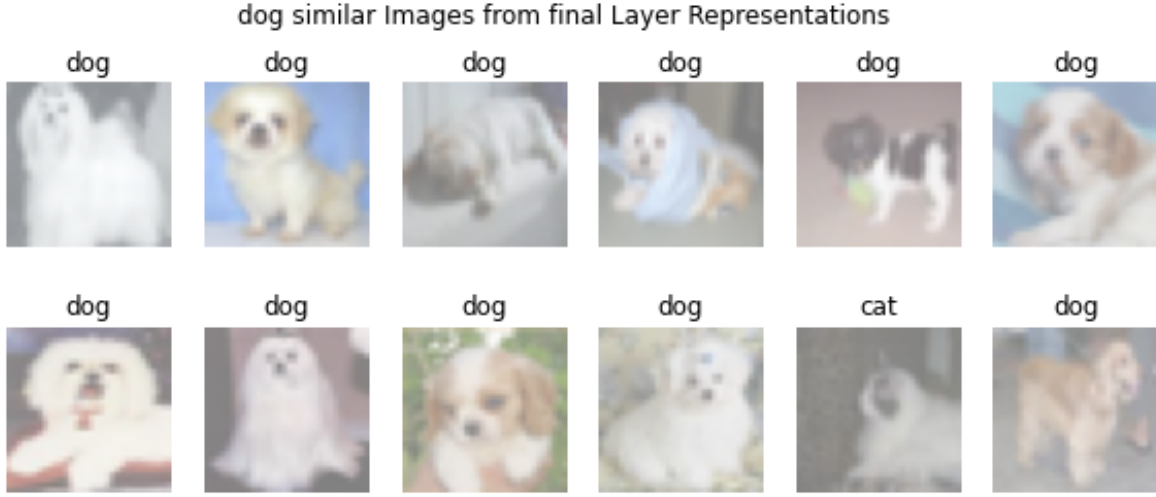
(a) Plane - End of Fully Connected Layer Representation



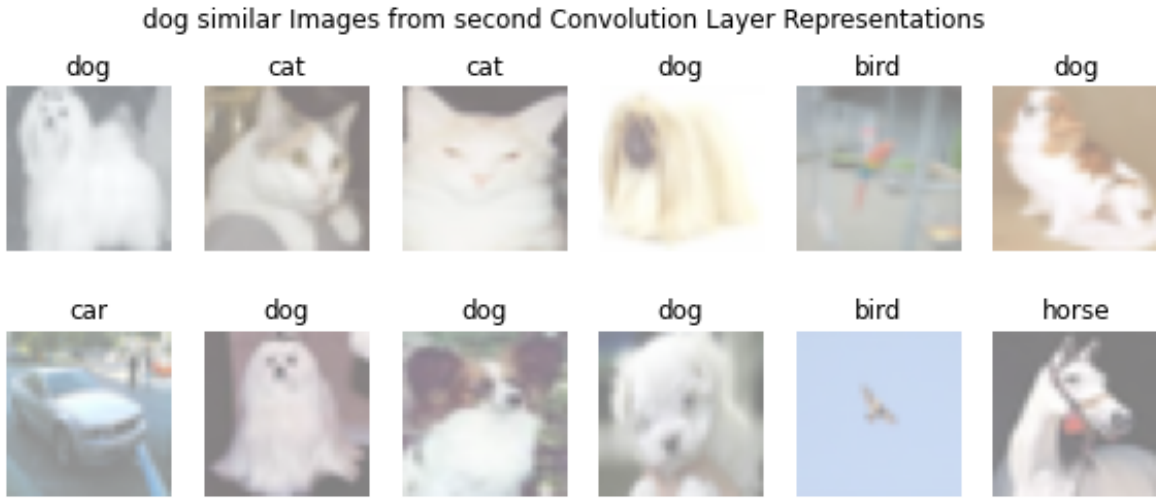
(b) Plane - End of Second Convolution Layer Block Representation

- **[DOGS SIMILAR IMAGES]** Taking another example for Dogs (*figure - 2 (a) & (b)*), we can observe that for similar images retrieved from the second convolution Layer input, we can see a couple of cat images, a horse, a bird, and a car. Maybe the high-level feature is the face with two eyes and a muzzle. That is why even the horse image has a very high similarity at this stage in the network. // On the contrary, when the images from the final layer representations are compared, we can see all of them as dogs except one cat image. Model might be learning that dog has two eyes, a round face and NO protruding ears (which might be the differentiating factor between other images), and that is why horse image and cat images which had protruding ears and high similarity before have not appeared in the top similarity results at the end of this layer.

Figure 2: Similar Images for the "first Image of first row" [DOG] in each section



(a) Dog - End of Fully Connected Layer Representation

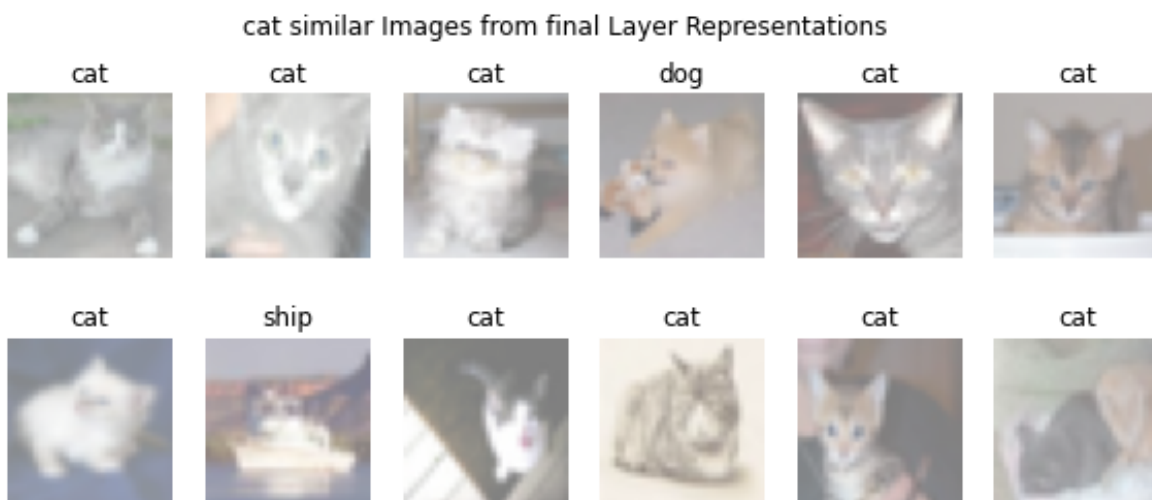


(b) Dog - End of Second Convolution Layer Block Representation

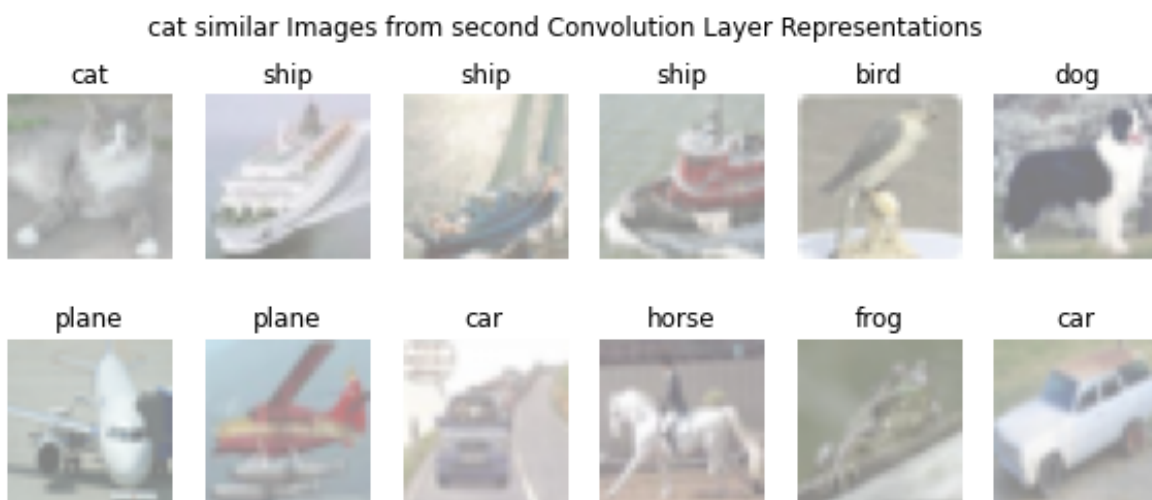
RESULT From the above examples representations, we can confidently say that the network starts with learning low-level features such as eyes, muzzle, noses, ears (whether protruding or not), etc and use these low-level features to map them to high-level features and use those to predict the final class.

Adding similar images for all the other classes. This gives more ideas on differentiating features between the classes.

Figure 3: Similar Images for the "first Image of first row" [CAT] in each section



(a) Cat - End of Fully Connected Layer Representation

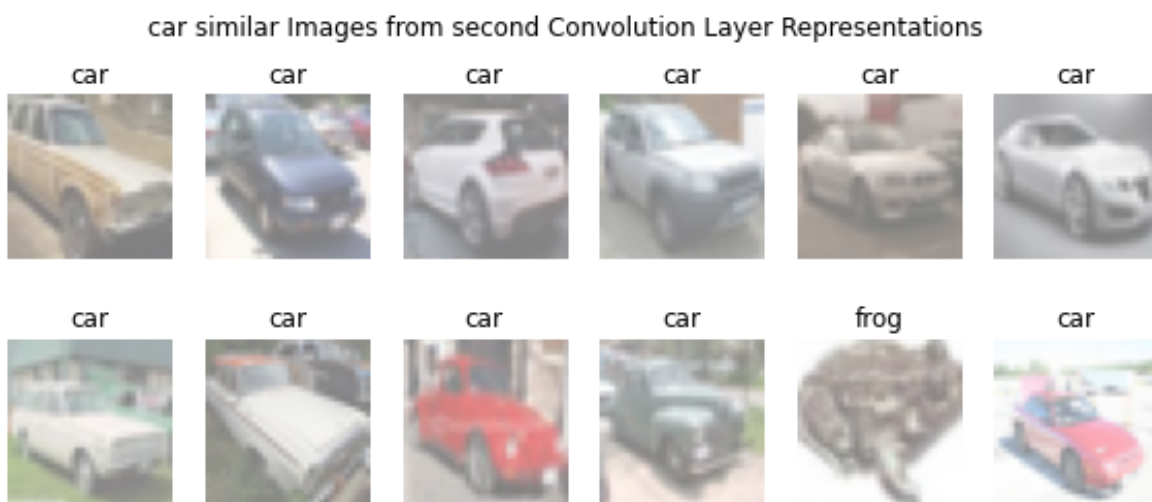


(b) Cat - End of Second Convolution Layer Block Representation

Figure 4: Similar Images for the "first Image of first row" [CAR] in each section

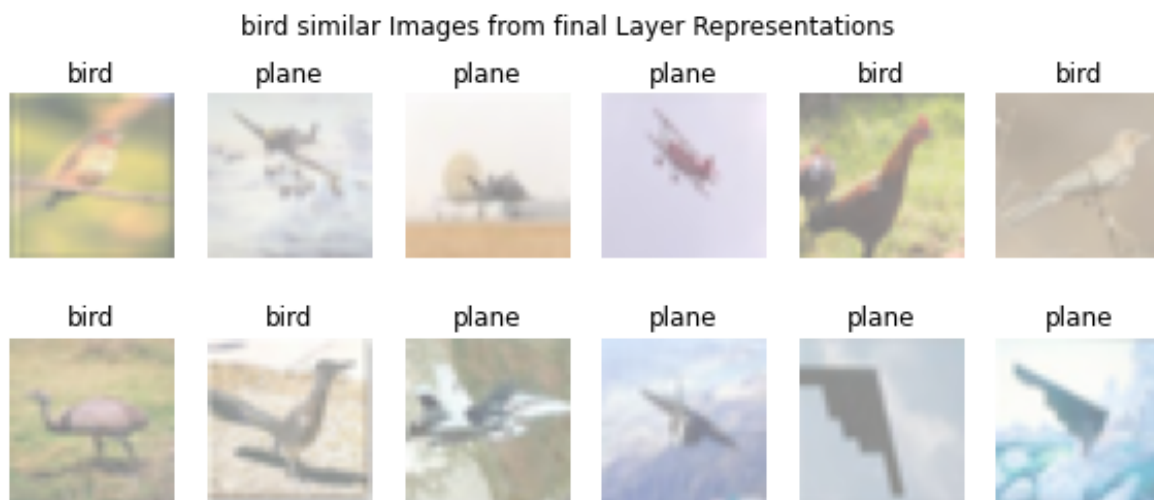


(a) Car - End of Fully Connected Layer Representation



(b) Car - End of Second Convolution Layer Block Representation

Figure 5: Similar Images for the "first Image of first row" [BIRD] in each section

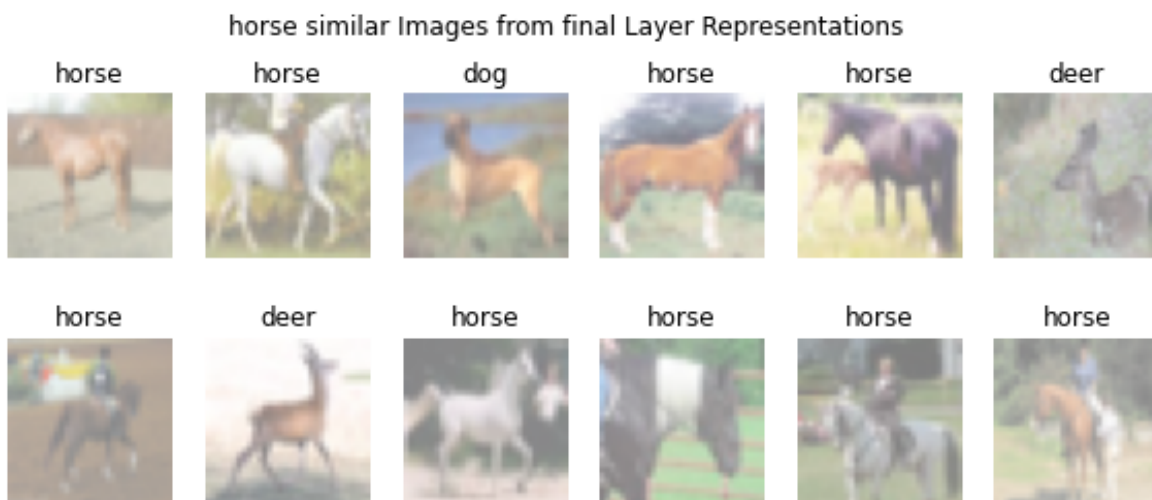


(a) Bird - End of Fully Connected Layer Representation

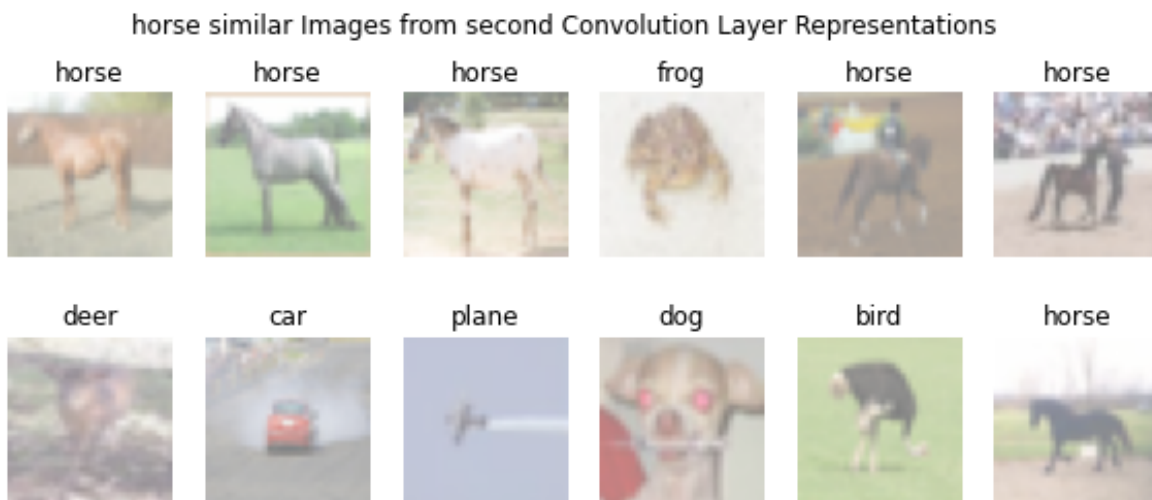


(b) Bird - End of Second Convolution Layer Block Representation

Figure 6: Similar Images for the "first Image of first row" [HORSE] in each section

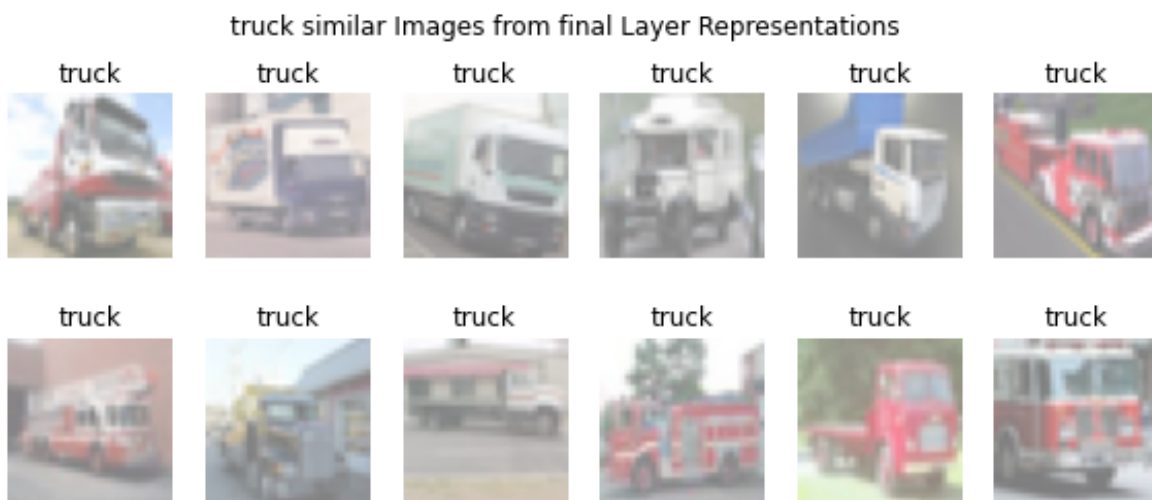


(a) Horse - End of Fully Connected Layer Representation

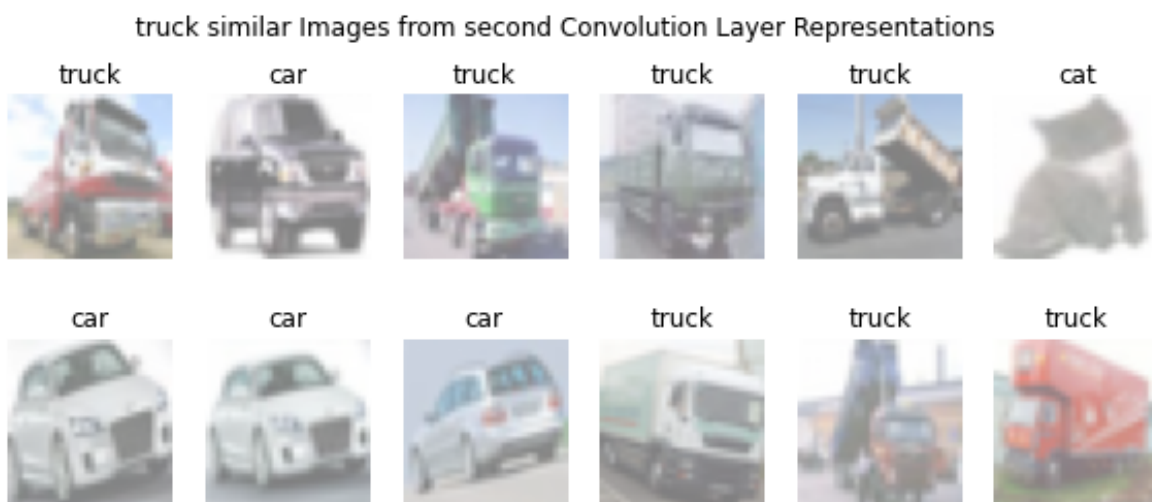


(b) Horse - End of Second Convolution Layer Block Representation

Figure 7: Similar Images for the "first Image of first row" [TRUCK] in each section



(a) Truck - End of Fully Connected Layer Representation



(b) Truck - End of Second Convolution Layer Block Representation

Figure 8: Similar Images for the "first Image of first row" [DEER] in each section

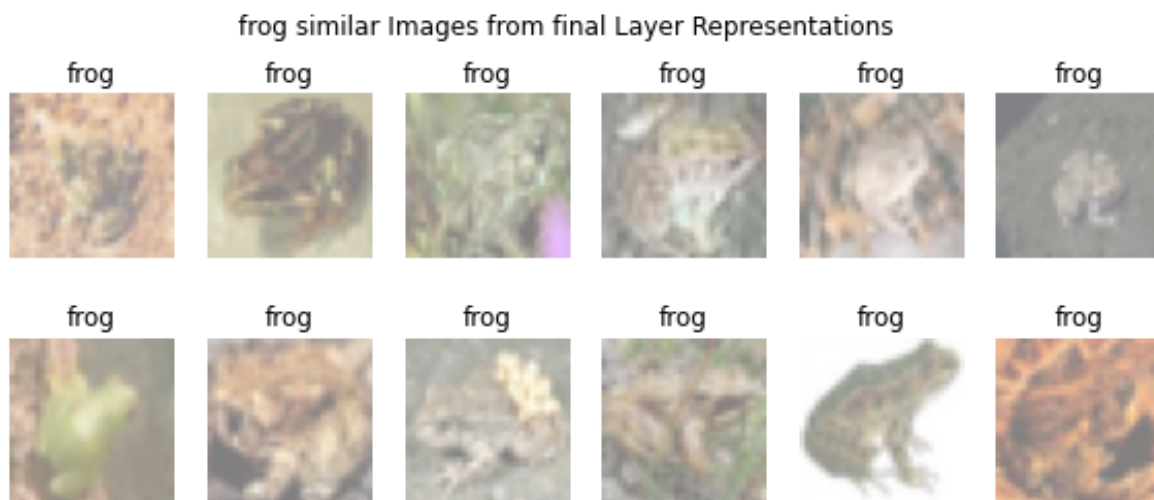


(a) Deer - End of Fully Connected Layer Representation

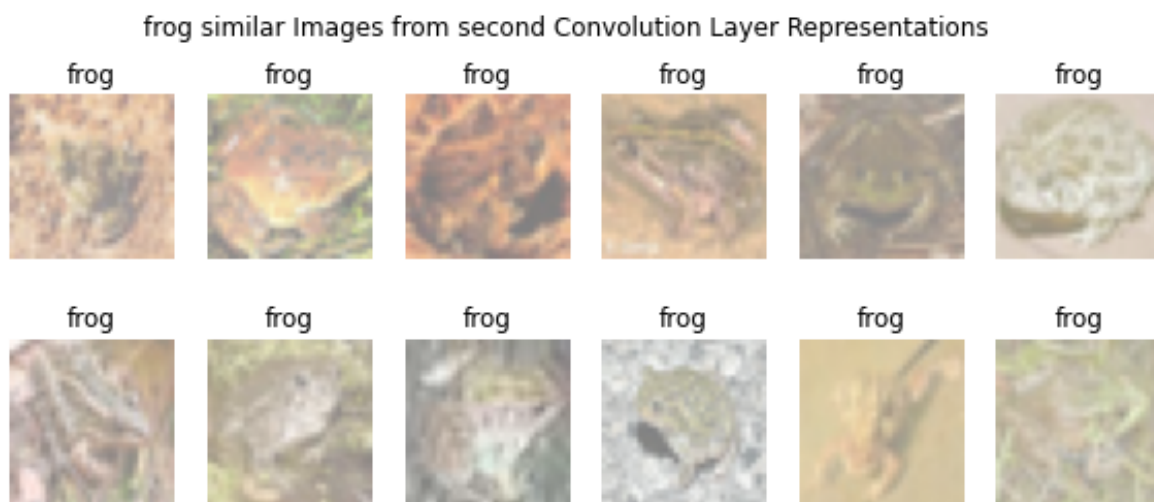


(b) Deer - End of Second Convolution Layer Block Representation

Figure 9: Similar Images for the "first Image of first row" [FROG] in each section



(a) Frog - End of Fully Connected Layer Representation

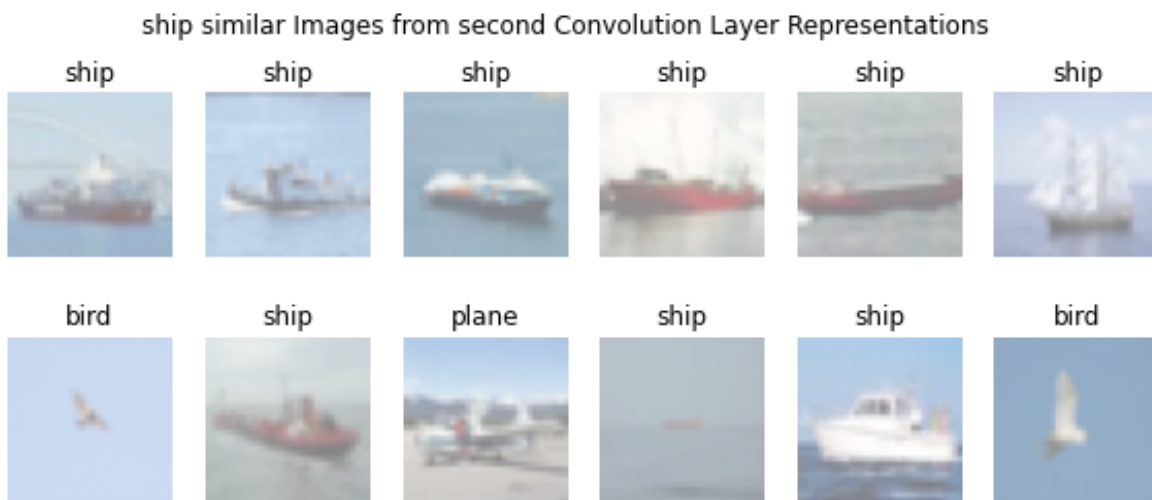


(b) Frog - End of Second Convolution Layer Block Representation

Figure 10: Similar Images for the "first Image of first row" [SHIP] in each section



(a) Ship - End of Fully Connected Layer Representation



(b) Ship - End of Second Convolution Layer Block Representation

3 CODE