
Exploring Transfer Learning, Fine-tuning of Thyroid Ultrasound Images

K.V. Sai Sundar, S. Siva Sankara Sai
Sri Sathya Sai Institute of Higher Learning
Prasanthi Nilayam, Puttaparthi, Anantapur-515134
saisundarkandarpa@gmail.com

Abstract

1 Deep Learning for medical imaging has been on the forefront of its numerous
2 applications, thanks to its versatility and robustness in deployment. In this paper
3 we explore various classification methodologies that are employed for datasets
4 of relatively small in size to actually train a deep learning algorithm from scratch.
5 Thyroid ultrasound images are classified using a small CNN from scratch, transfer
6 learning and fine-tuning of Inception-v3, VGG-16. We present a comparison of
7 the aforementioned methods through accuracy, sensitivity and specificity.

8 1 Introduction

9 1.1 Background

10 Image Classification task of a Machine Learning/Deep Learning algorithm is widely found to be
11 useful in the medical imaging domain wherein the algorithm is trained to classify medical images
12 into benign or malignant in case of cancer detection and various other ailments based on the symp-
13 toms in the images. Arriving at automated diagnosis seems to be the dream of every researcher
14 associated with computer vision coupled with biomedical imaging. Contrary to famous datasets
15 like the Imagenet dataset, Cifar-10, Cifar-100 which have thousands of images in each category,
16 the availability of medical data is limited. When running a deep learning algorithm with millions
17 of parameters, less data hinders the training with overfitting. The model eventually tends to fail at
18 generalization of learning giving low accuracy on the test dataset. Regularization can reduce the
19 high variance to some extent but training a deep learning framework from scratch remains out of
20 bounds. Data augmentation can be employed to further boost the dataset size. Therefore in order
21 to arrive at accuracies which can be of deployment standard, we resort to training a smaller CNN
22 from scratch, transfer learning-using bottleneck features from deep CNNs to train a new FC layer or
23 a different classifier and finally fine-tuning deep architectures to classify the custom dataset. We also
24 explored deployment of Thyroid Cancer Classification by training the mobilenet-224 with Thyroid
25 images and integrating a mobile application to do the classification. The code in this work is written
26 in TensorFlow Abadi et al. [2016] and Keras Chollet et al. [2015].

27 1.2 Motivation

28 The Thyroid Ultrasound domain was chosen following consultation with local doctors who brought
29 to our attention the high prevalence of thyroid ailments in the region. Thyroid Ultrasound is predom-
30 inantly used to detect thyroid nodules, classify them as benign or malignant and also to identify goi-
31 ter, thyroiditis. The problem consists of binary classification initially identifying images of patients
32 who probably require a biopsy in order confirm malignancy of nodule and eventually multi-class
33 classification identifying various other ailments apart from cancer. The aim of this work is to de-

34 velop an automated thyroid diagnosis system that could aid the radiologist and fasten the diagnosis.
35 In this paper we only discuss our implementations of binary classification.

36 **2 Dataset Description**

37 We have used two datasets in this work. The first dataset is a publically available one consisting
38 of 298 images and their corresponding biopsy verified reports in .xml format. Pedraza et al. [2015]
39 The TIRADS scores are given for each of the images ranging from 2 to 5 on the scale of increasing
40 probability for malignancy. Since our task in this work dealt only with probably benign or malignant
41 test scenario we considered scores 2 and 3 as benign and all the scores above these as malignant.
42 The second dataset used in this work was the local database of images from GE LOGIQ P9 which
43 were labelled by an experienced doctor and the reports were written in word format. The various
44 cases of cancerous nodule, thyroiditis, simple goiter, multinodular goiter, toxic goiter and normal
45 were present in this database. Again we considered only the relevant images as mentioned for the
46 previous dataset. This dataset consisted of thyroid images of 127 patients.

47 **3 Classification Methodology**

48 **3.1 Training a small CNN from scratch**

49 The first method constituted training a CNN from scratch using the medical data. The layer archite-
50 cture is 3 Convolutional Layers with 3x3 kernels of numbers 32, 32 and 64 respectively. The features
51 from the FC layer were classified using a regular sigmoid function in-to the two classes benign or
52 malignant. The training of the CNN was done on the GPU and since the model is relatively shallow
53 and the dataset small, the training completed in an hour's time.

54 **3.2 Transfer Learning**

55 Bottleneck features and the new FC Layer - The bottleneck features from the VGG-16 and Inception-
56 v3 were obtained and then trained on the CNN from the previous method. In VGG-16 the last three
57 FC layers were discarded and the CNN model which we used in the first method was fed these
58 features. The Imagenet pre-trained weights were loaded into the models and after the forward pass
59 of the image through the network bottleneck features were saved.

60 Bottleneck features and SVM- The CNN was replaced with the popular linear classifier Support
61 Vector Machine which was fed the bottleneck features for classification. The simple default parame-
62 ters of the SVM implementation provided by the scikit-learn were used in this meth-od. Deep CNNs
63 are known to be excellent feature extractors and using linear classifier to use these features proves
64 to be an excellent way of tackling smaller datasets. Razavian et al. [2014]

65 **3.3 Fine-tuning Inception-v3 and Vgg-16**

66 The Inception-v3 model was imported with the help of tf-slim high level API pro-vided by Tensor-
67 Flow. With the help of checkpoints provided for each of the models, pre-trained models could be
68 availed and fine-tuned. The Inception-v3 net provided in the slim API returns the list 'end points'
69 and 'logits' which can be fed to a classifier to predict the class. We obtained the end points['pre-
70 logits'] which is a layer prior to the last layer in the architecture and customized the FC layer, to
71 give output as a binary classifier. Softmax classifier was used for the classification. The last three
72 FC layers of the VGG-16 which contribute to huge computations were discarded and new FC layer
73 was attached after the 'pool-5' layer. For fine-tuning all layers above the conv 5_2 were frozen. So
74 essentially the last three layers (excluding the FC layers) and the new custom FC layer were trained
75 in this approach.

76 **4 Results**

77 The metrics used for evaluating the aforementioned methods are accuracy i.e. ratio of the total
78 number of correct predictions to the total number of images predicted, sensitivity, which gave an
79 indication of true positive rate and specificity for true negative rate. The classification was done

80 on both the datasets separately and on the combined dataset. It was observed that the first public
 81 dataset gave high sensitivity and low specificity while the second dataset gave just the opposite. This
 82 is due to the nature of the data wherein the first public dataset consisted of biased data with number
 83 of cancerous samples on the higher side. The local dataset on the other hand had the bias towards
 84 normal samples. This problem was handled by combing the datasets and also data augmentation
 85 achieved by flipping, rotating and adding noise to the existing images. The table below summarizes
 86 the results obtained on the combined dataset alone. This combined dataset consisted of 2525 training
 87 samples and 613 test samples. The metrics are tabulated as percentages.

SI No	Model	Accuracy	Sens.	Spec.
1	CNN from Scratch	0.82	0.89	0.82
2	VGG_16-Bottleneck+SVM	0.99	1	0.985
3	Inc_v3-Bottleneck+SVM	0.967	0.985	0.95
4	VGG_16-Bottleneck+New FC	0.94	0.96	0.92
5	Inc_v3-Bottleneck+New FC	0.98	0.99	0.97
6	VGG_16 Finetuning	0.76	0.84	0.68
7	Inc_v3 Finetuning	0.79	0.8	0.77

Table 1: Summary of Classification Results on Combined Dataset

88 5 Discussion

89 The delicate balance between the various evaluation metrics is important for performance analysis of
 90 deep learning algorithm in biomedical imaging. Owing to the fact that the dataset size is small Fine-
 91 tuning the Deep Architectures resulted was bettered by the other approaches. The amalgamation
 92 of linear classifier like the SVM to the architecture (especially the VGG-16) turned out to be the
 93 best model with stable metrics. Linear classifiers are considered to be useful when the deep neural
 94 networks is trained on datasets which are very different from domain in question. In order to take

95 References

- 96 Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu
 97 Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-
 98 scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016.
- 99 François Chollet et al. Keras, 2015.
- 100 Lina Pedraza, Carlos Vargas, Fabian Narvaez, Oscar Duran, Emma Munoz, and Eduardo Romero.
 101 An open access thyroid ultrasound image database. In *10th International Symposium on Medi-
 102 cal Information Processing and Analysis*, volume 9287, page 92870W. International Society for
 103 Optics and Photonics, 2015.
- 104 Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-
 105 the-shelf: an astounding baseline for recognition. In *Computer Vision and Pattern Recognition
 106 Workshops (CVPRW), 2014 IEEE Conference on*, pages 512–519. IEEE, 2014.