

1



Hochschule
Bonn-Rhein-Sieg
University of Applied Sciences



2

Master's Thesis

3

Sonar Patch Matching via Deep Learning

4

5

Arka Mallick

6

Submitted to Hochschule Bonn-Rhein-Sieg,

7

Department of Computer Science

8

in partial fulfilment of the requirements for the degree

9

of Master of Science in Autonomous Systems

10

Supervised by

11

Prof. Dr. Paul G. Plöger

12

Prof. Dr. Gerhard K. Kraetzschmar

13

Dr. Matias Valdenegro-Toro

14

January 2019

¹⁵ I, the undersigned below, declare that this work has not previously been submitted
¹⁶ to this or any other university and that it is, unless otherwise stated, entirely my
¹⁷ own work.

¹⁸ _____
Date

Arka Mallick

20 Accurately modelling features manually for sonar image matching has always
21 been very challenging. Only recently in Valdenegro et al. [41], instead of hand
22 designed features, a Convolutional Neural Network (CNN) was encoded to learn
23 the general similarity function to be able to predict a pair of sonar patches belong
24 to the different instance of same object or a different one. The result has been a
25 significant improvement over the state of the art methods. However the results need
26 to be further improved. The error in its prediction affects the overall performance of
27 object recognition, tracking etc high level tasks which are based on patch matching.
28 In this work more advanced CNNs are evaluated with the goal of improving the state
29 of the art results. Using DenseNet it was possible to predict binary classification
30 score of matching and non-matching cases with 0.955 average AUC of ten trials,
31 where the network was trained from scratch. This was a clear improvement over
32 the state of the art [41] results of AUC 0.894 on the same unseen test dataset. For
33 Siamese network with DenseNet branches the best performance was 0.921 average
34 AUC (10 trials). A Siamese network with VGG branches were also evaluated along
35 with Contrastive loss. This network's best performance was 0.949 average AUC (10
36 trials). To ensure the evaluation is effective thorough hyperparameter optimization
37 has been performed for all three methods.

38

Dedicated to the memory of my maternal uncle,

39

Subhas Chandra Ghosh,

40

without whose support my dream of pursuing higher

41

studies would not have been possible.

42

May he rest in peace.

Acknowledgements

My deepest thanks goes to my supervisor Dr. Matias Valdenegro-Toro. For several reasons. From technical aspect, his guidance and suggestions were invaluable and always spot on. At the beginning I spent several months training for deep learning and exploring other prerequisites for this work. He has been very supportive and motivating throughout the period of this thesis and even much before. In spite of very busy schedule, he spent a lot of hours helping me, many times even during weekends. It was a great experience working with him. Also, it was a long time aspiration to work in the domain of underwater robotics. A special thanks to him for giving me this opportunity.

Special thanks Prof. Dr. Paul G. Plöger for making this project possible through his supervision and invaluable inputs, also for the neural networks lesson, which was pivotal for deep learning training. I sincerely thank to Prof. Dr. Gerhard K. Kraetzschmar for kindly agreeing to supervise this project. Along with all the technical lessons, I also thank him for the wonderful opportunity of being able to take part in Robocup competitions, which were great experience for me to learn hands on programming and be in touch with great minds. In the context, my sincere thanks Santosh, Oscar and Shehzad who were great mentors and taught me many important lessons. I also thank Prof. Dr. Rudolf Berrendorf and the entire cluster team for providing such great computational environment for all the students. Without cluster I could have never finished this work.

I must thank my best friends Swarit, Pranika, Meemansa, Aaqib, Nour and Arshia for making my stay in Bonn stress-free and enjoyable. I also thank Adhideb and Abhilash, without their help I would not dare dreaming about masters in Germany. Finally, I thank my family and my partner Arshia for being supportive and encouraging through out and for being there always.

70	1 Introduction	1
71	1.1 Motivation	2
72	1.2 Problem Statement	3
73	1.2.1 Objectives	4
74	2 State of the Art	5
75	2.1 Convolutional Neural Network	5
76	2.2 Siamese Network	6
77	2.3 Sonar image matching techniques	7
78	2.4 Learning similarity function	8
79	3 Methodology	11
80	3.1 Data	11
81	3.1.1 Sonar image capture procedure	11
82	3.1.2 Data generation for deep learning	12
83	3.2 Architectures	14
84	3.2.1 Densenet two channel network	14
85	3.2.2 Densenet Siamese network	17
86	3.2.3 Contrastive loss	19
87	3.2.4 VGG network	20
88	4 Evaluation	21
89	4.1 Implementation and measurements.	21
90	4.1.1 Search for best hyperparameters	21
91	4.1.2 Grid search	21
92	4.1.3 Training process	22

93	5 Results	25
94	5.1 DenseNet-Siamese network	25
95	5.1.1 Hyperparameters to be evaluated	25
96	5.1.2 DenseNet growth rate and layers per block analysis	26
97	5.1.3 Total parameters analysis	34
98	5.1.4 Standard deviation across blocks	35
99	5.1.5 Finer grid search analysis	36
100	5.1.6 Number of filter analysis	39
101	5.1.7 DenseNet dropout probability analysis	40
102	5.1.8 Reduction and bottleneck analysis	40
103	5.1.9 Fully connected layer dropout analysis	42
104	5.1.10 Fully connected layer output size analysis	43
105	5.1.11 Batch size analysis	43
106	5.1.12 Learning rate and optimizer analysis	44
107	5.1.13 Final grid search	49
108	5.2 DenseNet Two-Channel	52
109	5.2.1 Hyperparameters to be evaluated	53
110	5.2.2 DenseNet growth rate and layers per block analysis	54
111	5.2.3 Pooling analysis	57
112	5.2.4 Basic network structure	58
113	5.2.5 Number of filter analysis	59
114	5.2.6 Reduction and bottleneck analysis	60
115	5.2.7 Total parameters analysis	60
116	5.2.8 Dropout probability analysis	61
117	5.2.9 Batch size analysis	62
118	5.2.10 Learning rate and optimizer analysis	62
119	5.2.11 DTC final grid search	63
120	5.3 VGG-Siamese network with Contrastive loss	67
121	5.3.1 Hyperparameters to be evaluated	67
122	5.3.2 Basic network structure	68
123	5.3.3 Flipped labels	68
124	5.3.4 Conv filters analysis	68
125	5.3.5 Kernel size analysis	69

126	5.3.6	FC units size analysis	69
127	5.3.7	Initializer	70
128	5.3.8	Batch normalization analysis	74
129	5.3.9	Dropout probability analysis	75
130	5.3.10	Learning rate and optimizer	76
131	5.3.11	Batch size analysis	77
132	5.3.12	Total parameter analysis	78
133	5.3.13	Final grid search	78
134	5.4	Comparative analysis	82
135	5.4.1	AUC comparison	82
136	5.4.2	Monte Carlo Dropout analysis	84
137	5.4.3	Real prediction analysis	86
138	5.4.4	Run time analysis	86
139	5.4.5	Ensemble	86
140	6	Conclusions	91
141	6.1	Contributions	91
142	6.2	Lessons learned	92
143	6.3	Future work	92
144		References	93

List of Figures

146	1.1	Use of convolutional network for learning general similarity function	
147		for image patches. The patches in the image are samples taken from	
148		the data used in this thesis. Inspired from Zagoruyko et al.[47] . . .	3
149	2.1	Figure is taken from S. Emberton et al. [12]. Light gets absorbed	
150		and scattered by the particles and objects in the water, which causes	
151		poor contrast and spectral distortion in sonar images.	7
152	2.2	Two channel network (left) and Siamese (right) CNN architectures	
153		used in Valdenegro et al., though the figure and inspiration is from	
154		Zagoruyko et al.[47]	8
155	3.1	Algorithm for matching and non matching pair of patch generation	
156		from Valdenegro et al [41]	12
157	3.2	Data processing pipe line	12
158	3.3	Some sample of the matching and non-matching pairs from Valdene-	
159		gro et al [41]	14
160	3.4	Example DenseNet architecture from [17].	15
161	3.5	Siamese Densenet structure.	17
162	3.6	VGG Siamese network with contrastive loss	20
163	4.1	Hyper parameters word cloud.	22
164	5.1	DenseNet layers per block and growth rate mean AUC analysis . . .	30
165	5.2	DenseNet layers per block and growth rate max AUC analysis . . .	31
166	5.3	Best growth rate analysis.	33
167	5.4	Cumulative growth rate analysis (histograms)	34
168	5.5	Total number of parameters analysis	34

169	5.6	Average standard deviation on AUC across networks with same	
170		number of dense blocks, e.g., 2-2 and 3-3 has 2 dense blocks, 2-2-2,	
171		3-4-5 has 3 dense blocks etc.	35
172	5.7	Finer search with top 10 architectures	36
173	5.8	Box and whisker plot representation of AUC predictions of 2-2-3	
174		and 2-2 architectures. First boxplot is for 2-2-3.	38
175	5.9	Number of filter analysis for different growth rates	39
176	5.10	DenseNet dropout probability analysis.	40
177	5.11	Evaluation of reduction and bottleneck and mean AUC(across 10	
178		trials). The max AUC obtained by one of the trial is also displayed	
179		separately for experiment with bottleneck and without bottleneck. .	41
180	5.12	The total parameters are varying with reduction. It also gets affected	
181		by use of bottleneck, how ever that is very minimal.	41
182	5.13	Hyperparameters in DenseNet-Siamese architecture	42
183	5.14	Optimal dropout (d) probabilities for FC layer analysis	42
184	5.15	Optimal FC filter size (f) analysis.	43
185	5.16	DenseNet-Siamese optimal batch size analysis.	44
186	5.17	In order, (a)Adam, (b)Nadam, (c)Adamax, (d)RMSprop, (e)Adadelta	
187		learning rate analysis	46
188	5.18	Comparison of different optimizers across different learning rates. .	47
189	5.19	Adadelta evaluation for high learning rates.	48
190	5.20	Boxplot visualization of top five DenseNet-Siamese network	50
191	5.21	Top five network statistical comparison in bubble plot	51
192	5.22	Area under curve for best Config 1	53
193	5.23	DenseNet two-channel architecture analysis	56
194	5.24	DenseNet two-channel avg pooling analysis	57
195	5.25	Average vs flatten pooling	57
196	5.26	DTC Number of filter size analysis	59
197	5.27	Reduction and bottleneck analysis	60
198	5.28	DTC Dropout probability analysis	62
199	5.29	Batch size analysis	62
200	5.30	Learning rate analysis	63
201	5.31	Overall best result for DenseNet two-channel 0.972 AUC (Single run)	65

202	5.32	Boxplot analysis for top configurations, with statistical outlook on	
203		10 AUC scores obtained for each cases.	66
204	5.33	Conv filters size analysis for VGG branches.	69
205	5.34	FC units and layers analysis for VGG branches.	70
206	5.35	Performance comparison of different combinations of Conv and FC	
207		initializers	72
208	5.36	Evaluation of how batch normalization affects prediction with in-	
209		crease in batch size.	75
210	5.37	Dropouts probability analysis for VGG branches.	76
211	5.38	Search results of best optimizer and learning rate	77
212	5.39	Batch size analysis	78
213	5.40	Final grid search results for VGG-Siamese network with Contrastive	
214		loss.	79
215	5.41	ROC for the highest auc recorded for the network with best configu-	
216		ration.	80
217	5.42	Further statistical analysis of top 4 configurations are visualized in	
218		box and whisker plots.	81
219	5.43	20 pairs of sonar patches, from the test dataset with index 1000 to	
220		1019, and corresponding mean prediction and standard of deviation	
221		out of 20 trials with Monte Carlo dropouts.	88
222	5.44	MC. predictions with highest std.	89
223	5.45	MC. predictions with lowest std.	89

List of Tables

225	1.1	Best result from Valdenegro et. al. 'Different' represents the under-	
226		lying test objects were different from the train objects. The binary	
227		prediction scores are presented here from the original results in [41]	3
228	4.1	Custom grid search space example.	22
229	5.1	Fixed hyperparameter values for the evaluation setup.	28
230	5.2	Finer search with 10 architectures	37
231	5.3	The search space for learning rate and optimizer best hyperparam-	
232		eters. Since both are related they need to be evaluated together. . .	45
233	5.4	The learning rate for which the optimizers recorded it's best mean	
234		AUC.	48
235	5.5	Optimizer Adadelata evaluation with high learning rates. Displayed	
236		results are after filtering the outlier cases.	48
237	5.6	Best hyperparameter values for final evaluation.	49
238	5.7	Best hyperparameter values for the DenseNet-Siamese	50
239	5.8	Fixed hyperparameter values for the evaluation setup.	55
240	5.9	The average and flatten pooling comparison results.	58
241	5.10	DenseNet two-channel 'basic network structure' for further hyperpa-	
242		rameter search.	59
243	5.11	Total number of parameters are compared across different number	
244		of dense block and different growth rates.	61
245	5.12	Best hyperparameter values for final evaluation of DenseNet two-	
246		channel.	63
247	5.13	Final grid search results	64
248	5.14	Best hyperparameter values for the DenseNet two-channel	66
249	5.15	The search space for VGG-Siamese network with Contrastive loss. .	67

250	5.16 Basic network configuration for the hyperparameter search. In this	
251	section of report this configuration will be referred to by the alias of	
252	'basic network structure'.	68
253	5.17 Kernel initializer top results	73
254	5.18 Batch normalization analysis for VGG branches.	74
255	5.19 The search space for learning rates specific for each optimizer for	
256	VGG-Siamese and Contrastive loss network.	77
257	5.20 Final grid search space VGG-Siamese network.	79
258	5.21 Best hyperparameter values for the Vgg-Siamese with contrastive	
259	loss	80
260	5.22 Comparative analysis on the AUC and total number of parameters	
261	in the best performing networks.	82
262	5.23 MC dropout prediction example for VGG-Siamese with Contrastive	
263	loss for a test image (index 1002). Compared for three models	
264	trained with different dropout rates.	85
265	5.24 Ensemble predictions on test data.	86

Introduction

“You can’t cross the sea merely by standing and staring at the water.”

- Rabindranath Tagore

More than two-thirds of the earth surface is covered by ocean and other water bodies. For a human it is often impossible to be able to explore it extensively. For finding new source of energy, monitoring tsunami, global warming or may be just to learn about deep-sea eco-system or may be to look for a lost ship or an airplane, the need for venturing into potentially dangerous underwater scenarios appear regularly, in fact getting more frequent. That’s why more and more robots are being deployed in underwater scenarios and lot of research is going in this direction. Some of the exploration or monitoring tasks requires the robot to see underwater, to make intelligent decisions. But underwater environment is very difficult for optical cameras, as light is attenuated and absorbed by the particles in the water. And lot of real life monitoring and mapping tasks take place in cluttered and turbid underwater scenario. The limited visibility range (underwater) of optical sensor is a big challenge. Hence sonar is more practical choice for underwater sensing as the sound can travel long distances with comparatively little attenuation.

A underwater robot, equipped with sonar image sensors, regularly needs to perform basic tasks such as object detection and recognition, navigation, manipulation etc. For most of it, effective acoustic vision is fundamental. In underwater

scenarios, sonar patch matching functionality is very useful in various applications such as data association in simultaneous localization and mapping (SLAM), object tracking, sonar image mosaicing [18] etc. Patch matching, in general, is heavily used in computer vision and image processing applications for low-level tasks like image stitching [3], deriving structure from motion [31], also in high-level tasks such as object instance recognition [28], object classification [46], multi-view reconstruction [37], image-retrieval etc. Typical challenges in patch matching tasks are different viewing points, variations in illumination of the scene, occlusion and different sensor settings. For sonar patch matching the common challenges with acoustic vision adds to the overall complexity. For example, low signal-to-noise ratio, lower resolution, unwanted reflections, less visibility etc. Because of these challenges the underlying object features might not be so prominent as a normal optical image. It has also been found that it is very challenging to manually design features for sonar images and popular hand designed features such as SIFT [29] are not always very effective in sonar images [41]. **That is why patch matching for sonar images remains a topic of further research interest.** – The ending needs to be improved?

1.1 Motivation

In Valdenegro et al. [41] it was first displayed that deep learning methods can be directly applied to the sonar image patch classification task. Without using any hand designed features it is still possible to perform a patch matching task with high accuracy. In [41] the authors recorded Forward-Looking Sonar (FLS) images of different objects in underwater environment and generated balanced dataset of image patches for the classification task. Two channel convolutional network and Siamese convolutional network were used in predicting the binary classification scores (table 1.1) on unseen test data. Using Two channel CNN the overall best result obtained was of Receiver operating characteristic (ROC) area under curve (AUC) of 0.894, with test accuracy of 82.9%. The reported results are better than both classic key-point matching methods (AUC 0.61 to 0.68) and machine learning based methods such as Support Vector Machine (SVM) and Random Forests (AUC 0.65 to 0.80). These findings were pivotal for the work presented in this thesis.

Network Type	Output	Test Objects	AUC	Mean Accuracy
2-Chan CNN	Score	Different	0.894	82.9%
Siamese CNN	Score	Different	0.826	77.0%

Table 1.1: Best result from Valdenegro et. al. 'Different' represents the underlying test objects were different from the train objects. The binary prediction scores are presented here from the original results in [41]

1.2 Problem Statement

However the results in [41] is far from perfect yet and the mis-classifications affect the object detection, recognition and tracking, which are dependent on the patch matching task. It is desired that the result is further improved. With this goal the same dataset has been obtained from the [41] for further evaluation. In this thesis more advanced architectures (such as DenseNet) [17] will be evaluated on the data and the goal is to improve the state of the art on the dataset.

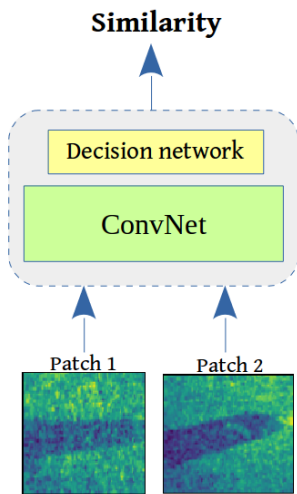


Figure 1.1: Use of convolutional network for learning general similarity function for image patches. The patches in the image are samples taken from the data used in this thesis. Inspired from Zagoruyko et al.[47]

For performing the patch matching task without using any hand designed features such as SIFT [29], we need to encode a network to learn the general similarity function (figure 1.1) for sonar image patches. Which enables the network to predict that two input sonar patches contain different views of a same object or not i.e matching or non-matching respectively.

With the goal of obtaining the best possible result in the evaluation, best hyperparameters are needed to be found out. The hyperparameters are the parameters of network or learning process whose values are already set before the training starts [45]. For example the starting learning rate or the optimizer can be fixed before the training process starts.

343

344 1.2.1 Objectives

- 345 1. Select and evaluate a number of state of the art architectures which can
346 perform sonar patch matching more effectively and improve the state of the
347 art result (table 1.1). The dataset presented in [41] is used for evaluating
348 different architectures.
- 349 2. Through grid search, determining the best hyperparameters for each of the
350 architectures is very important part of this work.
- 351 3. Another objective is to present a comparative analysis of the performance
352 of all the architectures. The metric of comparison is the ROC area under
353 curve value of the prediction score computed on the test dataset (D) from
354 Valdenegro et al [41].

355 The structure of the report is as follows. In the State of the art chapter 2 the
356 related works that are done is presented. In Methodology chapter 3 the data and
357 architectures evaluated are presented in details. The Evaluation chapter contains
358 details of how the grid search was performed for finding best hyperparameters.
359 Then in Results chapter the best hyperparameters search results and eventually
360 best hyperparameters for each of the architecture and their best prediction on the
361 test data is presented. The Results section ends with comparative analysis of all
362 the architectures evaluated. Lastly in Conclusions, the contribution of this work,
363 future research directions etc. is discussed.

State of the Art

2.1 Convolutional Neural Network

Before starting with the related works, little introduction to the Convolutional neural network (CNN) is provided here. CNN [44] is a special type of Artificial Neural Network (ANN) with Multi-Layer Perceptrons (MLP), consisting of learnable weights and biases. CNNs are usually trained with back propagation algorithm in a supervised manner. While a CNN also expresses a single differentiable score function similar to a normal ANN, CNNs are different in architecture though, specially designed for recognizing visual patterns directly from the raw image pixels with minimal preprocessing [26]. Unlike fully-connected architecture, in CNN, neurons in a layer will only be connected to a small region of the layer before. CNNs take advantage of the explicit assumption that inputs are images, as a result they can be more efficiently encoded to have much lesser parameters. The main building blocks of CNNs are Convolutional Layer (Conv), Pooling Layer, Fully-Connected Layer (FC). Now a typical example of a CNN can be denoted by [Input - Conv - ReLU - MP - FC], where Input holds raw pixel values of the input image, Conv layer applies convolution operation to the local region of input and produce output; this operation is inspired by the natural response of a neuron in biological visual cortex getting stimulated. ReLU layer provides activation to each elements and applies a threshold ($\max(0, x)$) so that the negative weights become zero. Pooling layer which basically down-samples the input along the

spatial dimensions of supplied height and width, by mapping a cluster of neurons from one layer to a single output in the next layer. MP or max-Pooling layer takes the maximum value from the cluster of the neurons and maps in the next layer. **Average** pooling takes average values of the cluster and maps to the next layer. Pooling can also be local or global. Lastly the FC layer just connects every neurons of one layer to the next layer. Fully-Connected layers serve as the decision network and determines the overall output size.

Inspired by [41] following notations for CNN layers are used to describe components of the architectures: **Conv(Nf, Fw x Fh)** is a convolutional layer with Nf filters of width Fw and height Fh. A max-pooling layer is represented by **MP(Pw, Ph)** where Pw x Ph is the sub-sampling size of the layer, width and height respectively and **FC(n)** is a fully connected layers where n represents output size.

In the following section Siamese network is briefly discussed because it is important to have a understanding of it before some of the state of the art network architectures can be explained in details.

2.2 Siamese Network

In 1993 the concept of a new artificial neural network called Siamese network was introduced by Bromley et al. [2]. Siamese network consists two identical sub-networks which are joined at the output. During training stage one of the input pairs is connected to each of the subnetworks. The main idea here is that the subnetworks (also called branches) are trained simultaneously and extract features from the inputs. While the shared neurons are capable of measuring the distance between the two extracted feature vectors by each branch. If the predicted distance between two feature vectors is lesser than a threshold then it can be considered that the inputs are similar, otherwise non-similar. Authors used this concept to compare two signatures in [2]. One of the signature was previously obtained from the authentic owner (up to 6 signatures were recorded). This was then used to compare with a new signature to verify if the both persons are same or not, as precaution against forgery. Authors implemented the Siamese network to be able to extract and compare different features of the two signatures and if the output is within a threshold then they were considered matching. If not then it was most

likely a forgery.

2.3 Sonar image matching techniques

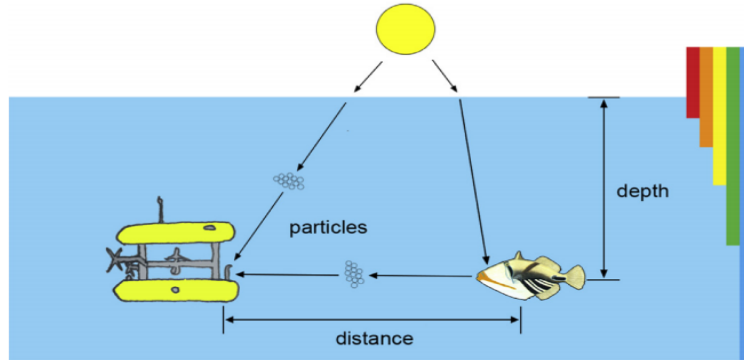


Figure 2.1: Figure is taken from S. Emberton et al. [12]. Light gets absorbed and scattered by the particles and objects in the water, which causes poor contrast and spectral distortion in sonar images.

Sonar image patch matching is more difficult than normal optical patch matching problem. This is because sonar images have additional challenges such as, non-uniform insonification, low signal-to-noise ratio, poor contrast [12], low resolution, low feature repeatability [19] etc. But sonar image matching has important applications like in sonar registration, mosaicing [24], [18] and mapping of seabed surface [32] etc. While Kim et al. [24] used Harris corner detection and matched keypoints to register sonar images, Hurtos et al. [18] incorporated Fourier-based features for registration of FLS images. S. Negahdaripour et al. [32] estimated mathematical models from the dynamics of object movements and its shadows. Vandrish et al. [42] used SIFT [29] for sidescan sonar image registration. Even though these approaches did achieve considerable success in respective goals, were found to be most effective when the rotation/translation between the frames of sonar images are comparatively smaller. Block-matching was performed on segmented sonar images by Pham et al. [33], using Self-Organizing Map for the registration and mosaicing task. Recently, [48] for stereo matching, [23] for fast under-water object detection and localization, [40] objectness scoring, CNNs are being more and more used for sonar image processing. The main reason behind such rise of CNN usage is that, the CNNs can learn sonar-specific information from the sonar data

directly. No complex manual feature design or rigorous data preprocessing steps are needed, which makes the task less complex but prediction accuracy achieved is actually higher.

2.4 Learning similarity function

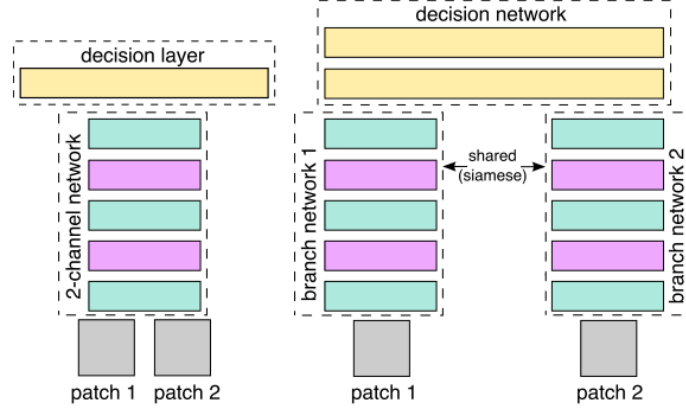


Figure 2.2: Two channel network (left) and Siamese (right) CNN architectures used in Valdenegro et al., though the figure and inspiration is from Zagoruyko et al.[47]

Zagoruyko et al. have demonstrated that CNNs can be directly deployed to learn the underlying similarity function to be able to determine that two input images/patches are different instances of same object or not, without any help from hand designed features. The scarcity of accurate hand designed features for sonar images motivated Valdenegro et al. to evaluate similar approach as Zagoruyko et al. and encode a CNN for sonar image comparison. Valdenegro et al. [41] evaluated two architectures, a two-channel network and a Siamese network. The architectures figure 2.2 were based on the work from Zagoruyko et al. [47]. A grid search was used over predefined set of hyperparameters to encode the best performing network. The final two channel network structure presented in the paper for predicting binary classification score was as follows, Conv(16, 5 x 5)-MP(2, 2)-Conv(32, 5 x 5)-MP(2, 2)-Conv(32, 5 x 5)-MP(2, 2)- Conv(16, 5 x 5)-MP(2, 2)- FC(64)-FC(32)-FC(1). Similarly grid search was also conducted for Siamese network structure for classification score. The structure for each of the branches or

sub-network is as follows, Conv(16, 5 x 5)-MP(2, 2)-Conv(32, 5 x 5)- MP(2, 2)-Conv(64, 5 x 5)-MP(2, 2)-Conv(32, 5 x 5)-MP(2, 2)-FC(96)-FC(96).

The output features from the branches were then concatenated to form 192 element vector. This was then passed through a decision network with FC(64)-FC(1). For both two channel and Siamese score prediction network, the final FC layer had Sigmoid [9] activation function and binary cross entropy loss [6] function.

As noted from [41] not very complex network structures were used, but the obtained results are very good. In theory, more advanced networks or loss functions which help extracting more discriminative and patch invariant features, should improve the results even further.

Methodology

3.1 Data

Just like any other machine learning project the data is the most important part of the evaluation. Hence for the proper qualitative evaluation a deeper look at the data and how it was generated is needed. The data set for training and testing are both obtained from Valdenegro et al [41]. Following is a brief description of the overall pipeline to obtain the data for training.

3.1.1 Sonar image capture procedure

In Valdenegro et al [41] the authors have explained in details how the raw sonar images were captured. In a water tank in their facility, the images were taken using ARIS Explorer 3000 and mounted Forward-Looking Sonar (FLS). The objects featured are household garbage items and common marine debris. There are total 9 different types, such as metal cans, bottles, drink carton, metal chain, propeller, tire, hook, valve. In [41], in controlled underwater environment total of 2072 images; about 2500 instances of the aforementioned objects were captured and labeled with 9 classes accordingly.

```

1:  $\bar{L}_m \leftarrow \emptyset, L_{nm} \leftarrow \emptyset$ 
2: for  $\text{img} \in I$  do
3:   for object  $o \in \text{img}$  do
4:      $OC \leftarrow \text{crop } B_o \text{ from img.}$ 
5:     for  $i = 0$  to  $S_p$  do
6:        $MC \leftarrow \text{sample random object } p \text{ of class } C_o \text{ and}$ 
        $\text{make an image crop.}$ 
7:       Append  $(OC, MC)$  to  $L_m$ 
8:     end for
9:     for  $i = 0$  to  $S_n$  do
10:       $NMC \leftarrow \text{sample random object } p \text{ of class } C_p \neq$ 
       $C_o, \text{ and make an image crop.}$ 
11:      Append  $(OC, NMC)$  to  $L_{nm}$ 
12:    end for
13:    for  $i = 0$  to  $S_n$  do
14:       $BC \leftarrow \text{sample random background patch and}$ 
       $\text{make an image crop.}$ 
15:      Append  $(OC, BC)$  to  $L_{nm}$ 
16:    end for
17:  end for
18: end for

```

Figure 3.1: Algorithm for matching and non matching pair of patch generation from Valdenegro et al [41]

484

3.1.2 Data generation for deep learning

485 As part of the same work [41], matching and non-matching pairs of sonar image
 486 patches were generated using a patch generation algorithm displayed in figure 3.1
 487 where each patch, an instance of one of the 9 classes, were obtained from meaningful
 488 crops of the original sonar images. Authors have figured that 96x96 pixels is the
 489 best size of the crops.
 490

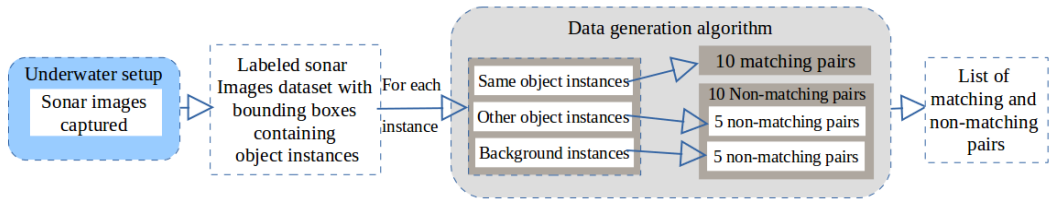


Figure 3.2: Data processing pipe line

491 To create matching pair, two patches that belong to the same object-type
 492 but different perspective and insonification level were used. For generating non-

matching pair, two instances of different object-types were used. Also, one instance of an object and a background patch, which contains no object instance, were used to generate additional non-matching patches. According to the authors of [41], for balanced and effective learning, the ratio of matching and non-matching pairs in both train and test data is maintained at 1:1. That is, for every ten matching pair five object-object non-matching and five object-background non-matching pairs were generated. In figure 3.2 the overall pipeline for the data generation is displayed in simple block diagram. In figure 3.3 some sample patches from all type of pairs are displayed. Using the patch generation algorithm total of 39840 pairs were generated from the instances of 6 distinct object type, which were used as the train data. While another 7440 pairs were generated from the instances of remaining object types, for testing purpose. This test dataset does not contain any common object with the training dataset, it should be a good test for the generalization of the approaches. The labels for the data are 0 and 1 representing non-matching and matching pairs respectively. For this thesis work, this dataset is obtained in HDF5 files [16]. It contains patches in form of tensors [39], each of shape (2,96,96). Here the pair of patches (size 96x96) are placed in a way that each channel contains one patch.

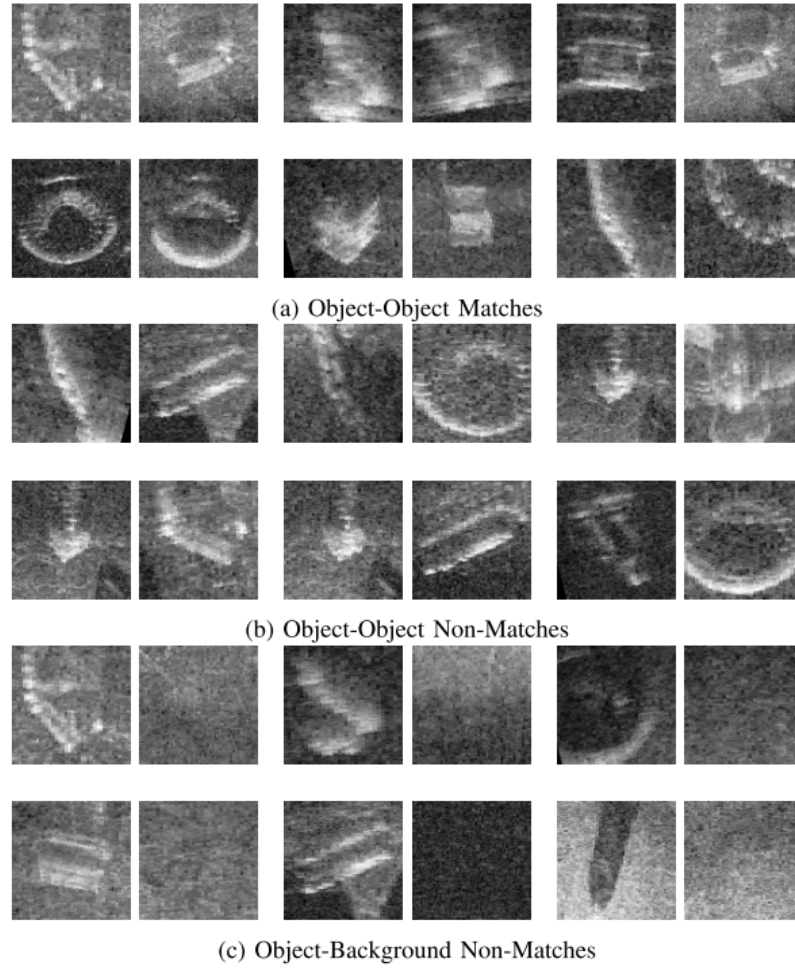


Figure 3.3: Some sample of the matching and non-matching pairs from Valdenegro et al [41]

511 3.2 Architectures

512 In this section the theoretical aspect of the different architectures that we are
 513 using are explained along with brief implementation details where applicable.

514

515 3.2.1 Densenet two channel network

516 In Densenet each layer connects to every layer in a feed-forward fashion. With
 517 the basic idea to enhance the feature propagation, each layer of Densenet blocks
 518 takes the feature-maps of the previous stages as input.

The Densenet implementation used is from the keras community contributions (keras_contrib) [30]. Another Densenet implementation was initially evaluated from Thibault de Boissiere [10], but it was older implementation and did not have bottleneck or compression implemented.

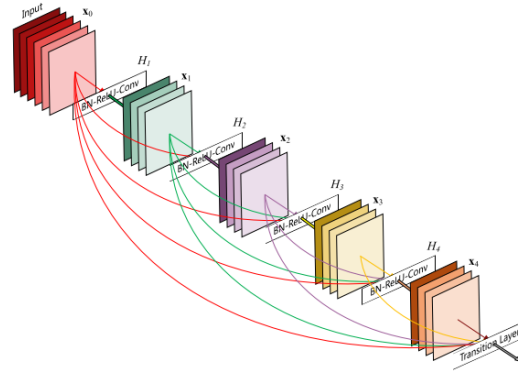


Figure 3.4: Example DenseNet architecture from [17].

In Densenet two channel the the sonar patches are supplied as inputs in two channels format, the network by itself divides each channel into one patch and learn the features from the patches and then finally compare them using the Sigmoid activation function at the end with FC layer of single output. The original label in the data is 0 for non-matching pair of input patches. Label 1 means matching pair.

Densenet architecture depends on the hyperparameter settings of the instantiation. The parameters and their brief definition is described in the following section.

- **Growth rate:** Number of filters to add per dense block. Growth rate regulates how much information is contributed by each layers to the global state. Global state is the collective knowledge of the network, that is the state of the previous layers are flown into each layer in form of feature-maps, which is considered as global state. Each layer adds k more feature-maps to the current state, when growth rate is k .
- **Nb_filter:** initial number of filters. -1 indicates initial number of filters will default to $2 * \text{growth_rate}$.

- 539 • **Nb_layers_per_block**: number of layers in each dense block. Can be a -1,
540 positive integer or a list. If -1, calculates nb_layer_per_block from the network
541 depth. If positive integer, a set number of layers per dense block. If list,
542 nb_layer is used as provided. Note that list size must be nb_dense_block.
- 543 • **Depth**: number of layers in the DenseNet.
- 544 • **Nb_dense_block**: number of dense blocks to add to end.
- 545 • **Bottleneck Layers**: To improve computation efficiency a bottleneck layer
546 with 1x1 convolution is introduced before each 3x3 convolution layers.
- 547 • **Compression**: Reduces the feature maps in transition layers and makes the
548 model more compact and computationally efficient.

3.2.2 Densenet Siamese network

In this architecture the branches of the Siamese network are Densenet. Following the classic Siamese model the branches share weights between them and get trained simultaneously on two input patches and then learn the features from the inputs. Through the shared neurons in the Siamese network it is able to learn the similarity function and be able to discriminate between the two input patches. The role of the Densenet branches is feature extraction, the decision making or prediction part is taken care of by the Siamese network.

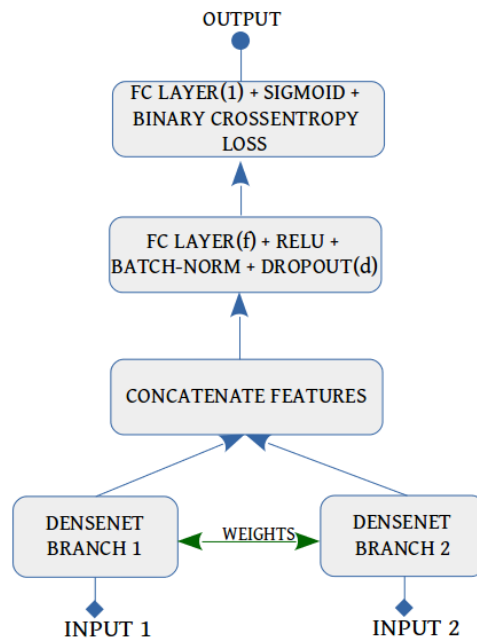


Figure 3.5: Siamese Densenet structure.

In figure 3.5 the basic architecture is displayed for the Densenet Siamese network. The two Densenet branches are designed to share weights between them. The extracted features are concatenated and connected through a FC layer, followed by activation Relu and where applicable Batch normalization and dropout layers. The output is then connected to another FC layer with single output, for binary prediction score of matching (1) or non-matching (0). Sigmoid activation function and binary cross entropy loss function is used for this final FC layer. As mentioned

565 in the figure 3.5 the size of the output of the FC layer (f) and value of dropout
566 probability etc. hyperparameters are to be decided after thorough hyperparameter
567 search.

3.2.3 Contrastive loss

Using contrastive loss higher dimensional input data (e.g. pair of images) can be mapped in the much lower dimensional output manifold, where similar pairs are placed closer to each other and the dissimilar pairs have bigger distances between them depending on their dissimilarity. So using this loss function the distance between two input patches projected in the output manifold can be predicted and if the distance is closer to 0 then the input pairs are matching, otherwise its dissimilar (above threshold). The formula for the loss is shown below.

$$D_w(\vec{X}_1, \vec{X}_2) = \|G_w(\vec{X}_1) - G_w(\vec{X}_2)\|_2$$

$$L(W, Y, \vec{X}_1, \vec{X}_2) = (1 - Y)\frac{1}{2}(D_w)^2 + (Y)\frac{1}{2}\{\max(0, m - D_w)\}^2$$

Here L is the loss term, the formula presented here is the most generalized form of the loss function, suitable for batch training. \vec{X}_1, \vec{X}_2 represents a pair of input image vectors. Y are the labels, 0 for similar pair and 1 for dissimilar pair. D_w is the parameterized distance function to be learned by the algorithm. m is the margin and m is always greater than zero and it defines a radius around G_w . The dissimilar pairs only contribute to the loss function if their distance is within the radius. One of the idea for evaluating this loss function is to use it with a Siamese network, as the loss function takes a pair of images as input. So its very relevant to the problem statement of this thesis.

To evaluate the contrastive loss a Siamese network has been selected since it also takes two input images and compare them. For the branch structure of the Siamese (figure 3.6) VGG network has been chosen. The role of the VGG network is to extract features, similar to the Densenet Siamese, the final decision making and prediction is taken care of by the Siamese network. In this case the output is actually euclidean distance between the two input sonar patches, projected into lower dimension using contrastive loss.

- Explain here about the Siamese network with VGG branches used and why Siamese network is chosen to be evaluated with the contrastive loss
- Do a comparative study with binary cross entropy too and see how that goes, probably it's difficult now, another way is to use this with the Siamese Densenet network and see.

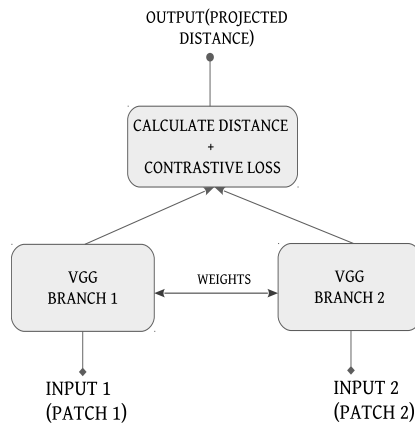


Figure 3.6: VGG Siamese network with contrastive loss

The basic VGG network structure is as follows, Conv(n , $a \times a$)-Conv(n , $a \times a$)-MP(2, 2)-Conv($2n$, $a \times a$)-Conv($2n$, $a \times a$)-MP(2, 2)-Conv($4n$, $a \times a$)-Conv($4n$, $a \times a$)-Conv($4n$, $a \times a$)-MP(2, 2)-Conv($8n$, $a \times a$)-Conv($8n$, $a \times a$)-MP(2, 2)-Conv($8n$, $a \times a$)-Conv($8n$, $a \times a$)-MP(2, 2)-Conv($8n$, $a \times a$)-Conv($8n$, $a \times a$)-MP(2, 2)-FC(d). The a , n and d are hyperparameters which needs to be optimized. The details of the evaluation is presented in the section 5.3

3.2.4 VGG network

Implementation taken from <https://hackernoon.com/learning-keras-by-implementing-vgg16-from-scratch-d036733f2d5>

Evaluation

4.1 Implementation and measurements.

- Compare based on AUC
- What is roc AUC
- AUC vs Accuracy, why AUC is better than Accuracy or other point based methods.
- Talk about tensors, tensorflow and keras here too
- Criteria of comparison, mainly the mean AUC then max AUC, then number of parameters and execution time*

4.1.1 Search for best hyperparameters

What are the hyperparameters. Effect of setting best hyper parameters.

4.1.2 Grid search

Currently the search space for the evaluation is hand designed and supplied to the evaluation script externally. Example



Figure 4.1: Hyper parameters word cloud.

Epochs	Optimizer	Batch size	Dropout	Use batch norm
5	adam	64	0.4	True
5	adam	64	0.4	False

Table 4.1: Custom grid search space example.

629 In this example it is displayed that how the search space is hand designed
 630 to evaluate different cases, here the effect of using batch normalization can be
 631 investigated by comparing the results of two test cases. All the values set in
 632 the grid passed as input to the actual code and defines the structure or different
 633 parameters from inside the code. In 4.1 the flag use batch norm is set to 'True',
 634 which gets propagated in the actual block of code where it enables the use of batch
 635 normalization layer, where applicable.

636 For the evaluation the epochs are set after multiple testings in order to ensure,
 637 to some extend, the networks does not get too overtrained and the generalization
 638 gets poorer. Also if the network is undertrained, then it will not generalize well.
 639 The setting the epochs were one of the big challenge of this work, because in some
 640 cases the networks gets overtrained after 4-5 epochs.

641

642 4.1.3 Training process

643 The train data has been divided into train and validation data randomly using
 644 sklearn train test data split [36] in 8:2 ratio. This is done using fixed seed to

ensure the validation and train cases are same for all the evaluations across all three methods. It is important to clarify that the roles of the validation data in keras is different than classic training method. In keras the validation data is not used to update the weights. Instead this validation dataset can be used to monitor the network performance in terms of validation accuracy and validation loss. Using callback functions such as `EarlyStopping` and `ReduceLROnPlateau` [4] the network validation performance can be monitored after each epochs of training. By configuring the callbacks properly (parameter: Early stop patience) the training can be stopped if the validation accuracy or loss does not improve after a number of epochs then the training process terminates. Similarly, if the validation performance does not improve after some epochs the learning rate can be reduced by a previously determined factor using `ReduceLROnPlateau` to come out of the local minima. The patience values are set after some manual trials.

```
#Early stopping
es = EarlyStopping(monitor='val_acc', patience=es_patience,verbose=1)
#Model check point
checkpointer = ModelCheckpoint(filepath=weight_file, verbose=2,
    save_best_only=True)
#Learning rate reducer on plateau
lr_reducer = ReduceLROnPlateau(monitor='val_loss', factor=np.sqrt(0.1),
    cooldown=0, patience=lr_patience, min_lr=0.5e-6,verbose=1)
```

Results

In this section the evaluation results of all three architectures and related best hyperparameter search results and detailed analysis is presented. Also the architectures are compared to each other based on the prediction on test data.

5.1 DenseNet-Siamese network

How the search for best hyperparameters for the DenseNet-Siamese are designed and corresponding results are presented in the following section.

5.1.1 Hyperparameters to be evaluated

The whole search space is divided into smaller blocks such as hyperparameters for DenseNet, hyperparameters for Siamese and others.

Hyperparameters of DenseNet

Layers per block defines both the depth of the DenseNet, which is automatically calculated and also defines the number of dense blocks. Besides growth rate, reduction, bottleneck, number of filters (nb_filter), pooling, include top, dropout, subsample initial block and weights are the most of the hyperparameters needed to be defined to instantiate a DenseNet. While some of this values will be fixed for all the evaluation, some needs to be chosen from a possible set of value, which could be infinite as well.

689 **Hyperparameters of Siamese Network**

690 Rest of the Siamese network that serve as decision network that connects the two
691 DenseNet branches have the hyperparameters as follows. FC layer output size,
692 dropout probability value.

693 **Other hyperparameters**

694 Apart from the aforementioned ones there are other general hyperparameters such
695 as learning rate, batch size for training and optimizer. In the following part of the
696 report it is described that how the whole search space is divided into smaller, well
697 defined parts. The results are visualized with help of charts and tables to help in
698 decision making.

699

700 **5.1.2 DenseNet growth rate and layers per block analysis**

701 Most important part of this network is the feature extraction capability of the
702 DenseNet branches. The overall performance of the network depends on it. So the
703 first focus of the hyperparameter search is to find out the main parameters defining
704 the DenseNet architecture. The overall search space for this could be very big or
705 even infinite. So for the first evaluation we define a coarse search space. With
706 hope to find a best performing parameter configuration or at least narrow down
707 the search space. Layers per block are chosen among 2, 3, 4, 5, 6. For single dense
708 block evaluation goes up-to 12 layers. More than that(14) causes memory to run
709 out as the size gets too big for the cluster gpu memory(16GB). Each network has
710 been evaluated for growth rates of 6, 12, 18, 24, 30, 36. Different dense block sizes
711 of 1, 2, 3, 4. The parameters compression/reduction and bottleneck are set 0.5 and
712 'True' respectively. Both this parameters control the compactness of the model
713 and help reducing the parameter required, hence in theory, making it possible
714 to evaluate much bigger networks without running into memory shortage issue.
715 The network that is being evaluated here are named DenseNet-BC by authors, i.e
716 DenseNet with bottleneck and compression.

717 There are other parameters but number of dense block, growth rate and layers
718 per block and reduction ratio are the main parameters which controls the architec-

ture and parameter size of the network the most. The goal of this focused search is to narrow down the overall search space from the DenseNet parameters perspective. For more fine grained analysis the compression and bottleneck parameters will also be evaluated.

DenseNet parameters number of filter (`nb_filter`) value are fixed at 16 for this search. The parameter classes are set to 2, which represents 1 for matching and 0 for not-matching pairs. 96, 96 is the input image dimensions. And it is single channel. So depending on local setting of the keras, 'channel-first' or 'channel-last' suitable input shape is chosen automatically as (1, 96, 96) or (96, 96, 1) respectively.

Subsample initial block enables the sub sampling of the initial input image to reduce the computation cost. The value is set to 'True'. The DenseNet parameter 'weights' value is set to 'None' to ensure that previously trained weights are not used. The decision network is not required for the branches as Siamese provides that, the value for parameter include top is set to 'False'. The learning rate used for the test is 2E-4. As a regularization measurement dropout probability value for DenseNet used as 0.2 to handle over-fitting. To ensure the grid search is effective, too much regularization is not good, as it can be restricting the overall evaluation at times. Epochs can be different for each architectures to ensure that the networks are able to achieve good training accuracy. But networks should not be over training, so the choice of epoch is selected after multiple manual trials for each set of network configuration in the search space. From the Siamese side of the parameters, after concatenating the DenseNet branch output feature maps, the combined features then passed through a fully-connected layer of 512 output size, which is followed by 'ReLU' activation and batch normalization (BN) and then a dropout layer with probability 0.5 has been added to ensure better generalization. Some of this values could have been further evaluated, how ever the values are obtained after lot of manual tuning and assured to be a decent starting point. 'Flatten' is used as pooling at the end of the DenseNet branches. Instead of the global average pooling from the original implementation. This causes increase in parameters overall though. Because with 'flatten' the multidimensional feature map at the end of DenseNet branch is just flattened. Where as in global average pooling [27], apart from the channel other dimensions are simply collapsed. But it is found that with 'flatten' the network is able to achieve much higher training

accuracy and generalization too for this network. Binary cross entropy loss function with 'Sigmoid' activation function used for the binary classification, this final layer acts as the binary classifier. In all the cases the networks are trained from scratch.

Growth rate and layers per block search setup summary

The overall search space is summarized in the section below.

Fixed hyperparameters

Other hyperparameters that are needed to instantiate the DenseNet are set to fixed values 5.1 after manual trials, in order to focus on the layers per block and growth rate parameters.

Table 5.1: Fixed hyperparameter values for the evaluation setup.

Name	Value	Name	Value	Name	Value
Number of filter	16	Subsample initial block	'True'	Weights	'None'
Dropout rate	0.2	Include top	'False'	Compression	0.5
Bottleneck	'True'	Pooling	'flatten'	Transition pooling	'max'
Siamese FC output	512	Siamese dropout	0.5	Optimizer	'adam'
Learning rate	2E-4				

Varying hyperparameters

- **Layers per block:**
 - **One dense block architecture (nb_dense_block=1)**
'2', '3', '4', '6', '8', '10', '12'
 - **Two dense block architecture (nb_dense_block=2)**
'2-2', '2-3', '2-4', '3-3', '3-4', '3-5', '4-4', '6-6'
 - **Three dense block architecture (nb_dense_block=3)**
'2-2-2', '2-2-3', '2-3-3', '2-2-4', '2-3-4', '3-3-2', '3-3-3', '3-3-4', '3-4-4',
'3-4-5', '3-3-6', '4-4-4', '4-4-2', '4-4-3', '4-4-6', '6-4-2', '6-6-3', '6-6-6'
 - **Four dense block architecture (nb_dense_block=4)**
'2-2-2-2', '3-3-3-3', '4-4-4-4', '6-6-6-6'

772 • **Growth rate:**

773 – Thin layers: 6, 12, 18

774 – Thick layers: 24, 30, 36

775 The evaluation result is thoroughly analyzed and presented in the following
776 section. Since the search space is big and each network configurations are evaluated
777 5 times, there are lot of data which are analyzed part by part with specific goals in
778 mind.

779 **Performance comparison based on mean AUC**

780 Each test case is trained from scratch and evaluated on test data 5 times. The metric
781 for evaluation is Area under curve(AUC) for the test data prediction. All together
782 there are too many results to display in report. So only **top 20** configurations with
783 highest mean AUC on test data across 5 trials are selected and displayed.

784
785 **Discussion**

786 From figure 5.1 it is observed that the DenseNet with two dense blocks or layers
787 (e.g., 2-2, 3-4) works best, ahead of single layer ones. Performance of the three and
788 four layer DenseNets are not good. This is unexpected. According to the original
789 paper [17] the performance of a normal DenseNet increases as more deeper the
790 network gets. However, in the original work the DenseNet is used individually as the
791 feature extraction and decision network both. In this case, DenseNet is only used
792 as the feature extraction network. Authors of [17] though hinted that the depth of
793 the network depends on the data volume available too. For example for ImageNet
794 [11], authors used four layer DenseNet with high growth rates. But they used three
795 layer DenseNet in other cases mostly. In any case when the DenseNet two-channel
796 network is evaluated this issue can be verified with more conviction than with
797 DenseNet-Siamese, simply because the DenseNet has been used in different way in
798 this case. From figure 5.1 no strong trend for growth rate is observed, so further
799 analysis or some other view of the data needs to be looked into. The network is
800 evaluated 5 times for each configurations because it is observed that not every

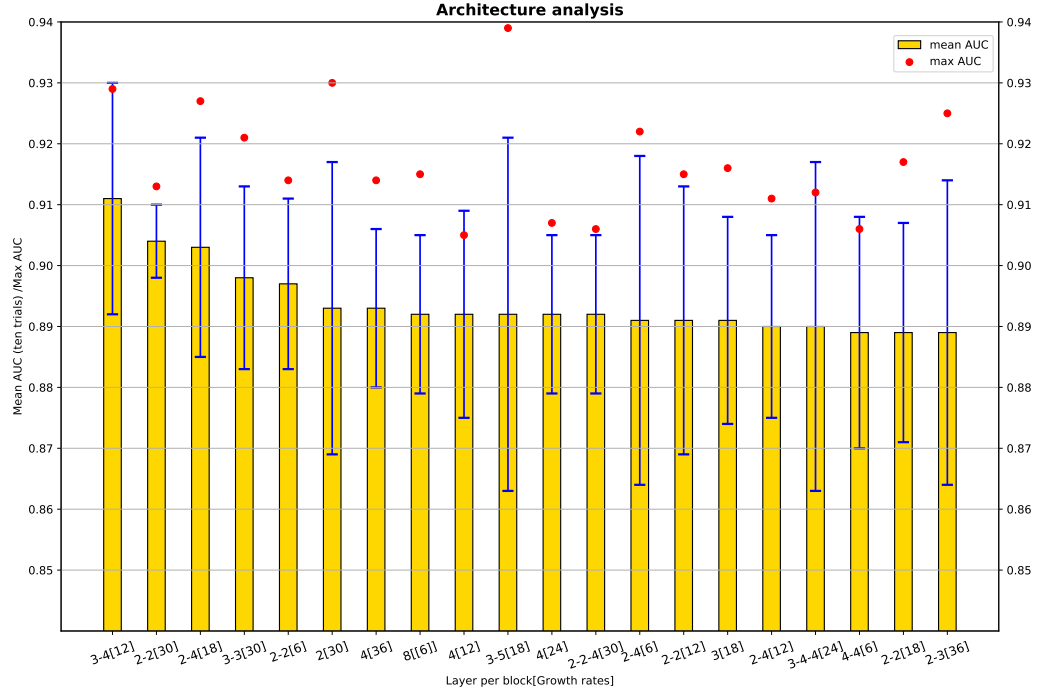


Figure 5.1: DenseNet layers per block and growth rate analysis, sorted by mean AUC, high to low. In x-axis the layers per block and growth rate are displayed together, with growth rate in brackets. It is easy to point out that, the top 20 is dominated by the two layers architectures. The four layer ones do not come close at all. Just few three layer and some single layer architectures are showing good results.

801 time the network performs exactly same way. Some times it score much higher and
802 some times it might even get stuck in a local minima. The weights of the network
803 are randomly initiated using default initializer Glorot uniform [7]. However many
804 parameters are involved and because every time the weights are drawn randomly
805 and the network gets trained from scratch, the decision boundary at the end of
806 same number of epochs may look very different. Since there are too many networks
807 to be evaluated only 5 times each of them are evaluated, with the idea that the
808 networks with good performance can be evaluated again for higher number of times
809 to verify their consistency.

Performance comparison based on maximum AUC

There are some networks, specially three layer ones, which has comparatively poorer mean AUC but at least one of the 5 runs they has scored very high AUC. That is why all the architectures are sorted according to their maximum AUC in one of the 5 trials. Displaying below in figure 5.2 is the **top 20** architectures (layers per dense block and growth rates) in terms of highest AUC on test data across 5 trials.

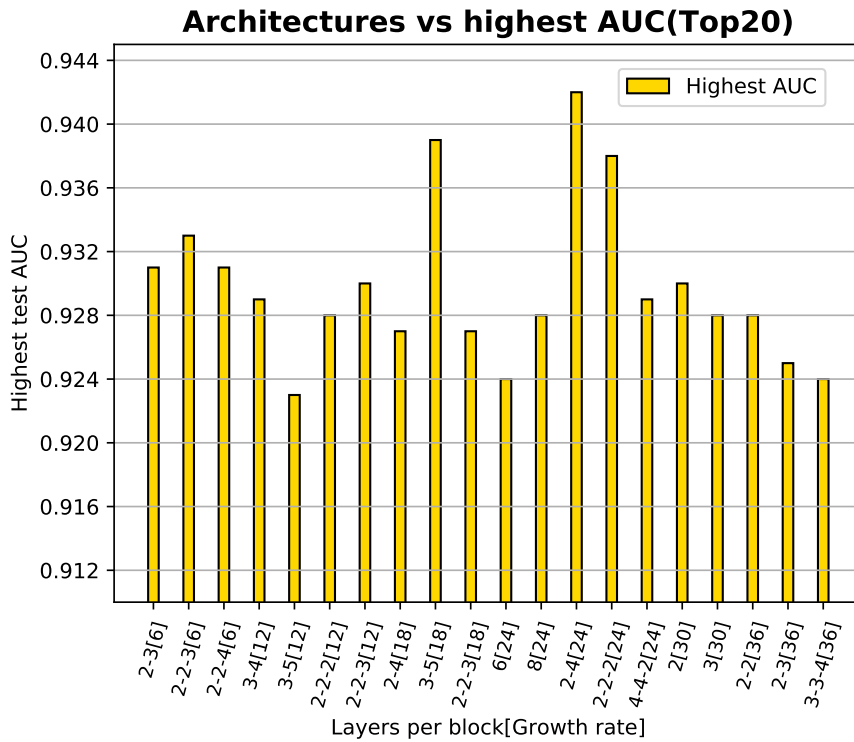


Figure 5.2: DenseNet layers per block and growth rate analysis, top 20 maximum AUC, sorted by growth rate. It seems the architectures growth rate 24 is most frequent here ahead of 12 and 18. unlike the top 20 mean AUC analysis many three layer architectures are in this list.

Discussion

In both top 20 lists mentioned above the four dense block architectures did not make it in any of the case. Their performance on the test data is worse, so under current setup and assumptions they work worse than lesser block networks.

Even though four layer networks are able to train above 95% train accuracy their generalization on the test data seems to be poor in general. Two dense block networks in general works best in terms of mean AUC of the 5 evaluations. In terms of max AUC score some of the Three block and one block DenseNet works very good as well, but may not be that consistent in general and did not make it in the top 20 mean AUC list. Some networks though, like 2-2, 2-4, 3-4, 3-5, 3, 8, 2-2-4 etc are of special interest since they have featured in both the list of mean and max top 20 AUC. Even though the top 20 maximum AUC list is dominated by architectures with growth rate 24, in top 20 architectures based on mean AUC that is not the case. So still unable to select the best growth rate for the further evaluation. Hence further analysis is done for the best growth rate on a different view of the data.

Optimal growth rate analysis

It is also important to find out the best growth rate for each of the architectures (layers per block). In the previous test each architectures are evaluated for growth rates: 6, 12, 18, 24, 30 and 36. For each architectures the growth rate for which the best mean AUC and best maximum AUC is recorded are displayed in the graph below (figure 5.3)

From figure 5.3 it is observed that the growth rate for which each architecture have best mean and for which it has maximum AUC, might not always be same. For this purpose, histogram of best performing growth rates obtained from mean and max AUC analysis is displayed in the figure 5.4 below:

Discussion

It is evident from figure 5.4b, based on the max AUC, there is no strong trend, its very random and inconclusive. Which is not very surprising given its just one run. With so many parameters and initialization and random dropouts involved, the network weights might learn very differently in spite of being trained in a same condition, resulting in a very different decision boundary. In figure 5.4a it is visible that the contribution of growth rate 6 and growth rate 36 is really less, so probably they are too thin or too thick for the data. From figure 5.4 above it is safe to assume that growth rate 18 is very good performer in both analysis. Which also

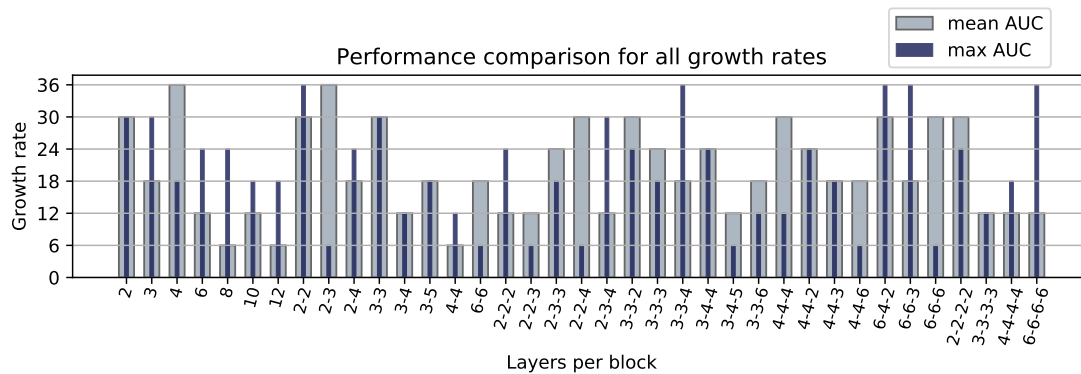


Figure 5.3: Best growth rate for each architectures are displayed based on mean AUC and max AUC. The grey column represents which growth rate ranked highest in mean AUC for the architecture, and the dark blue bar represents the growth rate for which the maximum AUC is recorded.

851 makes sense since it is neither too thin nor too thick. Because this conclusion is
 852 based on just 5 evaluations of each architectures, it make sense to experiment with
 853 other growth rates(except 6,36) as well for finer evaluation.

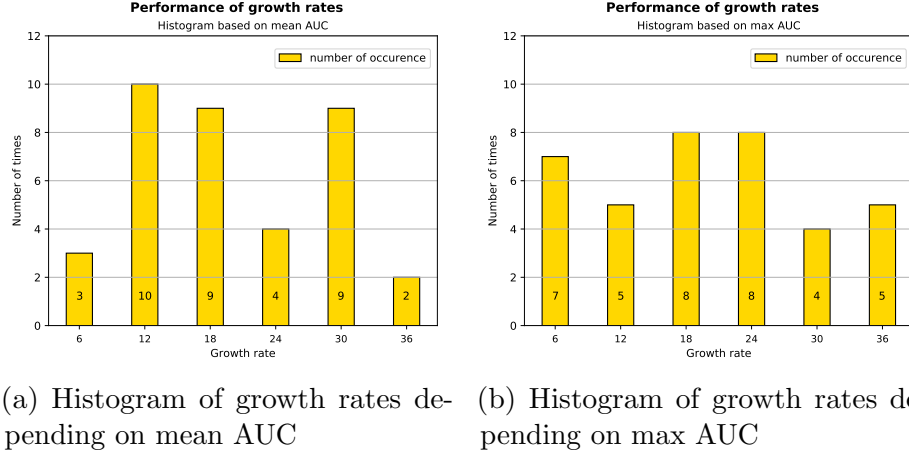


Figure 5.4: Cumulative growth rate analysis (histograms)

854

855 **5.1.3 Total parameters analysis**

856 It might be surprising but the single dense block networks have **most** parameters.
 857 This is because of the flatten pooling that is used here in this work instead of global
 858 average pooling 2D. And four block networks have the least total and trainable
 859 parameters. However as the number of dense blocks keeps getting higher the
 860 non-trainable parameters also gets higher. So 4 blocks dense net has most number
 861 of non-trainable parameters. For the visualization of the comparison 2, 2-2, 2-2-
 862 2, 2-2-2-2 layers per block DenseNet's parameter sizes are compared below, all
 863 recorded for growth rate 18.

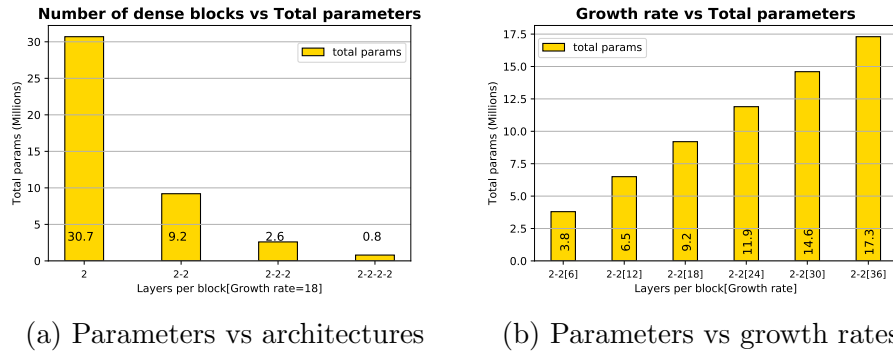


Figure 5.5: Total number of parameters analysis

In figure 5.5b it is shown that with growth rate increase the total parameter size also increases. For this purpose the parameters for all the growth rates compared for architecture 2-2.

5.1.4 Standard deviation across blocks

Another trend is observed that the standard deviation varies more as the number of dense blocks increase. So the standard deviation values for all the readings are collected for 1, 2, 3, 4 number of dense blocks (nb_dense_blocks) separately and their average values are presented in the table below.

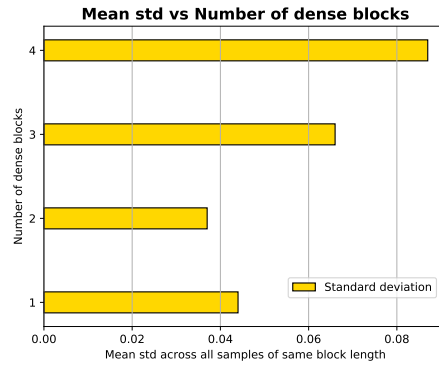


Figure 5.6: Average standard deviation on AUC across networks with same number of dense blocks, e.g., 2-2 and 3-3 has 2 dense blocks, 2-2-2, 3-4-5 has 3 dense blocks etc.

Even though the number of sample size differs a lot, 42, 48, 108, 24 samples for dense blocks 1, 2, 3, 4 respectively, it is observed that standard deviation of AUC is higher for dense block 3 and 4. It is probably because the blocks 3 and 4 networks have much lesser parameters than the 1 and 2 number of dense blocks. It is clear from this analysis that the two layer architecture is the best.

At this point, the best possible values of the growth rate and layers per block have been obtained. But the search space for the architectures are still big and needs to be shortened further and basically chose up to 3 networks so that further evaluations can be done.

883

884 **5.1.5 Finer grid search analysis**

885 From the first analysis the top 5 network based on mean AUC of 5 trials (figure
 886 5.1) and top 5 network obtaining maximum AUC (figure 5.2) are further evaluated
 887 for **20** evaluations each, after training from scratch. This 10 architectures will be
 888 referred as top 10 architectures in following analysis. All other test conditions remain
 889 the same. Results after 20 times evaluation is expected to be more dependable
 890 than 5 trials. In the chart 5.7 the mean AUC of 20 evaluations and it's standard
 891 deviation is displayed in yellow bars and blue lines respectively. While in red
 892 displayed the maximum AUC obtained in 20 evaluations.

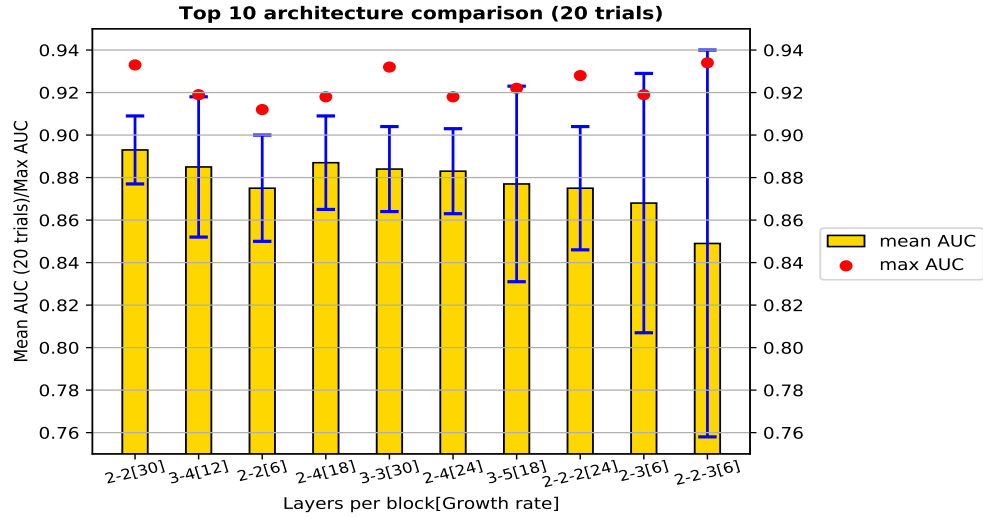


Figure 5.7: Finer search with top 10 architectures

893 Same data is displayed in table 5.2 but in decreasing order of the best mean
 894 AUC.

Table 5.2: Finer search results for 20 evaluation of previous top 10 architectures. Sorted according to mean AUC.

Layers	Growth rate	Mean AUC	Std	Max AUC
2-2	30	0.893	0.016	0.933
3-4	12	0.885	0.033	0.919
2-2	6	0.875	0.025	0.912
2-4	18	0.887	0.022	0.918
3-3	30	0.884	0.02	0.932
2-4	24	0.883	0.02	0.918
3-5	18	0.877	0.046	0.922
2-2-2	24	0.875	0.029	0.928
2-3	6	0.868	0.061	0.919
2-2-3	6	0.849	0.091	0.934

Discussion

It is observed that in 20 evaluations layers 2-2 with growth rate 30 is the best result both in terms of lowest standard of deviation and highest mean AUC. As it happens its also second highest in terms of the maximum AUC 0.933 just behind 0.934 from 2-2-3. 3-4, 2-4 networks are also performing well in terms of mean AUC. Their results are very close as well, so just evaluating 3-4 network for the finer analysis. 2-2-3 layers is interesting though, it has the highest standard deviation but 2 or 3 very good AUC scores too. So it needs to be further looked into.

In figure 5.8 the five number summary(min, first quartile Q1, median, third quartile Q3, max) is compared for architecture 2-2 and 2-2-3. For architecture 2-2-3, 3 AUC readings are detected as outliers out of 20 trials at 0.666, 0.608, 0.656 AUC. The outliers are affecting the overall mean auc for 2-2-3 network, also causing big standard deviation. This outliers are probably caused by training getting stuck in local minima or similar. 2-2 is found to be more consistent, it has no outliers. It is believed that consistency is desirable for a network. Hence 2-2-3 network is ruled out of contention for the best network, because of it's lack of consistency. though the overall concept of terming few prediction accuracies as outlier can be debatable. It is the nature of the network. At least, it is clear that three prediction accuracies are way of than other 17 predictions.

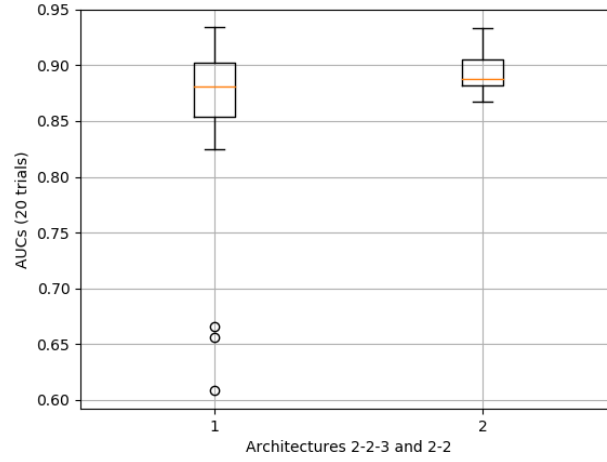


Figure 5.8: Box and whisker plot representation of AUC predictions of 2-2-3 and 2-2 architectures. First boxplot is for 2-2-3.

914 Five number summary

915 Also the brief introduction to the five number summary is as follows. **The**
 916 **minimum** is the smallest datum in the dataset. **The first quartile** is chosen so
 917 that 25% data points are lesser than it, similarly **the median** is the middle point
 918 (50%) of the dataset, and **the third quartile** is the point where 75% data falls
 919 below it. Lastly **the maximum** is simply the highest datum in the data set. All
 920 together this five criteria represents different aspects of any distribution.

921 Box and whisker plot interpretation

922 Box and whisker plot [13], where the box represent the interquartile range
 923 between first quartile 25% (Q1) and third quartile 75% (Q3) for the data range,
 924 and the orange line is for the median and the extended whiskers display the range
 925 of the maximum and minimum data points in the distribution. The interquartile
 926 range denoted by IQR, is the difference between Q3 and Q1. Outliers are the
 927 data points that reside outside the stretch of $[(Q1-1.5*IQR), (Q3+1.5*IQR)]$. The
 928 fraction 1.5 here is the default value in 'matplotlib' and also used in this work for
 929 all the evaluations. So if the data points lie outside the aforementioned range, then
 930 those are displayed as little dots or hollow circles in the plot beyond the whiskers.

5.1.6 Number of filter analysis

Initial number of filter (parameter name 'nb_filter') values 8, 16, 32, 64 are being evaluated here. Also a comparison between mean prediction AUCs obtained for growth rate 18 and 30 is done under this analysis. Ten evaluations of each test cases has been done.

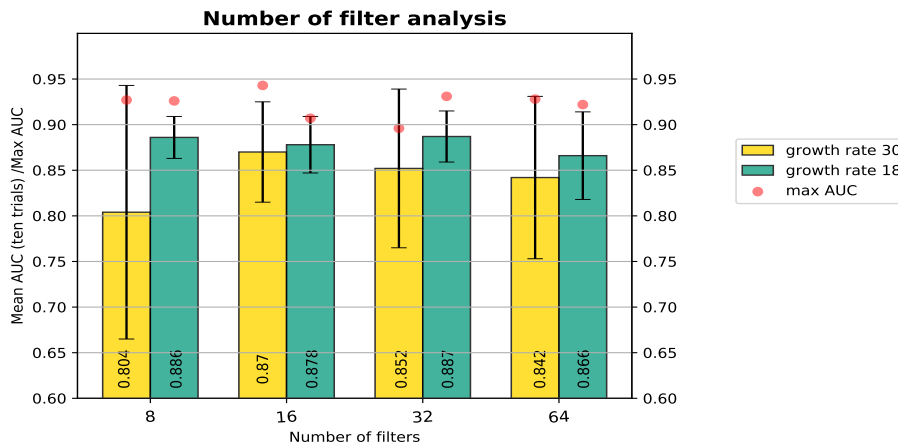


Figure 5.9: Number of filter analysis for different growth rates

Discussion

figure 5.9 shows that the number of filter 16 works better than others for growth rate 30 and also happen to have the Maximum AUC recorded and lowest standard deviation as well. It seems for thinner networks (lower growth rates), change in number of filter value matters lesser than the thicker networks. Since for 18 growth rate the mean AUC for all the values of 'nb_filter' are very close. But there is lot of difference for growth rate 30 for different rates. However both the networks are found to work pretty good with number of filter value of 16. This could have been further evaluated with more growth rates, but for this work it is concluded with 16 as our best nb_filter value.

947

948 5.1.7 DenseNet dropout probability analysis

949 As usual ten evaluations done for each dropout probabilities displayed in figure
 950 5.10. Mean AUC is best for dropout 0.4. Maximum AUC is highest for dropout 0.5.
 951 Now it is no surprise that with 0, 0.1 and 0.7 dropouts the results are not the best,
 952 because it's either too less or too much regularization. But it is bit unexpected
 953 to have dropout probability 0.3 and 0.5 performing low. probabilities 0.2, 0.4 are
 954 chosen as the best dropout rate values for future evaluations because they have
 955 comparatively better max AUC and mean AUC.

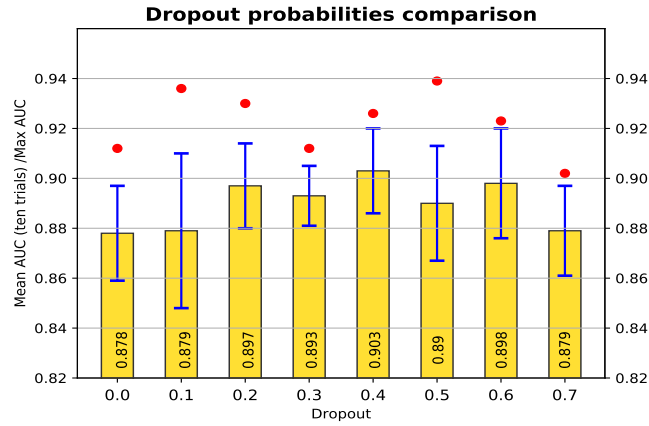


Figure 5.10: DenseNet dropout probability analysis.

956

957 5.1.8 Reduction and bottleneck analysis

958 The evaluation results for the bottleneck and reduction are displayed in the
 959 figure 5.11.

960 Discussion

961 The use of reduction really makes the model much more compact without losing
 962 the effectiveness. From figure 5.12 it is observed that without reduction the
 963 parameter size is 20.1 Millions, after using reduction of 0.7 the mean AUC is still
 964 as good but the total parameter size has become 12.1 millions. Use of bottleneck

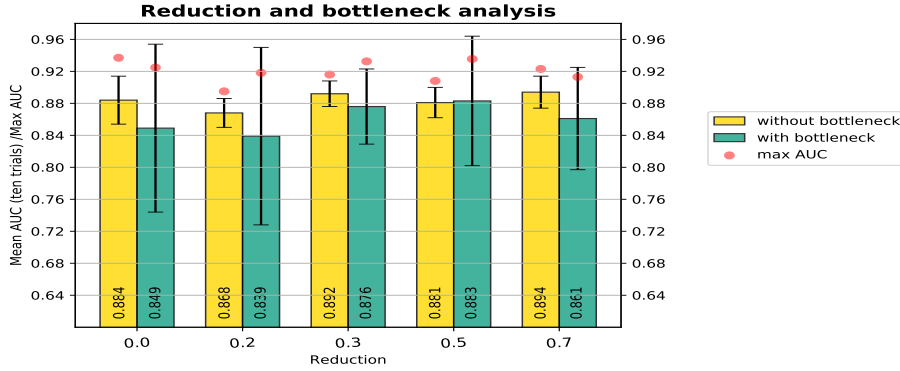


Figure 5.11: Evaluation of reduction and bottleneck and mean AUC(across 10 trials). The max AUC obtained by one of the trial is also displayed separately for experiment with bottleneck and without bottleneck.

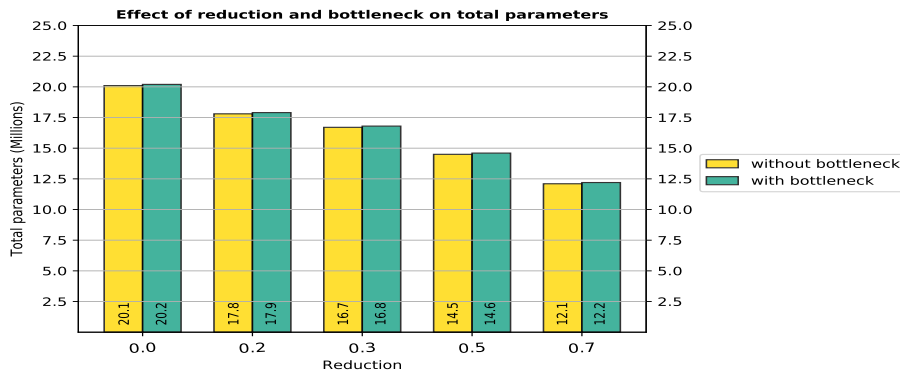


Figure 5.12: The total parameters are varying with reduction. It also gets affected by use of bottleneck, how ever that is very minimal.

965 does in-fact increase the parameters by little more than 0.1 million, but it does
 966 not improve the results at all. So this is a unexpected behavior. Bottleneck also
 967 has much higher standard deviation. Theoretically it is suppose to help making
 968 the model more compact. But over all effect on the data is observed to be adverse.
 969 But the max auc values still belong to the bottleneck layers 3 out of 5 times. Other
 970 two times also its close to highest. That's perhaps interesting to note. So the best
 971 reduction ratio chosen, based on the mean AUC are 0.7 and 0.3. Bottleneck will
 972 not be enabled for further evaluations.

973

974 **5.1.9 Fully connected layer dropout analysis**

975 In the diagram 5.13 the hyperparameters associated to the decision network
 976 part are displayed.

977 The dropout probability connected to the
 978 first FC layer of the Siamese part is evaluated
 979 here, shown as d in 5.13. The search space
 980 evaluated for the d is as follows: $[0, 0.1, 0.2,$
 981 $0.3, 0.4, 0.5, 0.6, 0.7]$.

982 **Discussion**

983 From figure 5.14 it can be observed that for
 984 dropout probability 0.5 the network has the
 985 highest mean AUC along with dropout proba-
 986 bility 0.7. The mean AUCs for 0.3, 0.4 is lower
 987 than expected and for 0.7 is much higher. The
 988 best values for dropout probability chosen for
 989 further evaluation are 0.5 and 0.7 both.

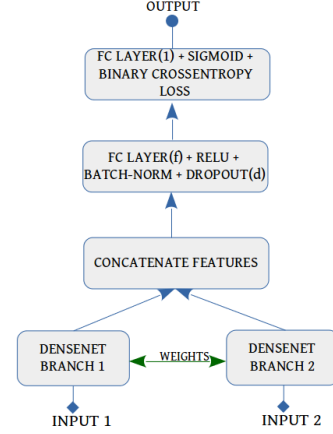


Figure 5.13: Hyperparameters in DenseNet-Siamese architecture

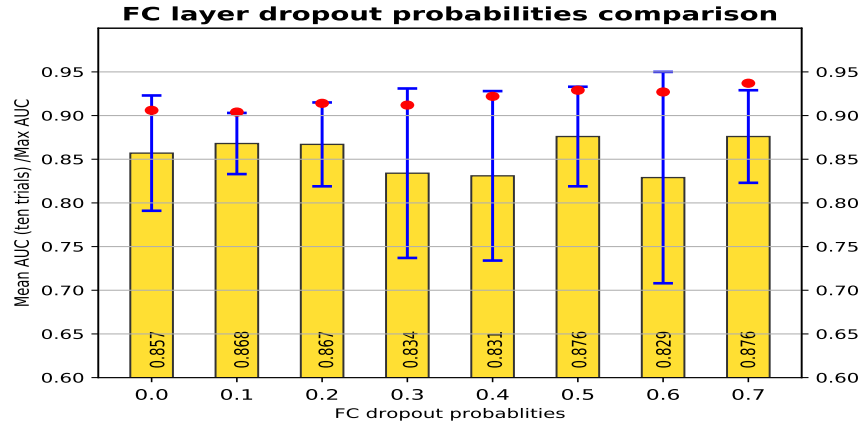


Figure 5.14: Optimal dropout (d) probabilities for FC layer analysis

990

5.1.10 Fully connected layer output size analysis

The FC layer output size is denoted by f in figure 5.13. The fully-connected layers are initialized (kernel) with He normal initialization.

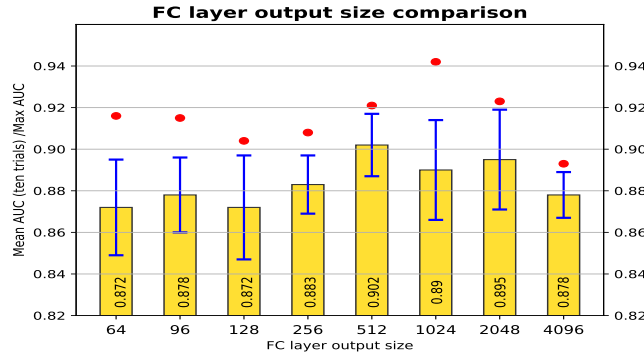


Figure 5.15: Optimal FC filter size (f) analysis.

Discussion

Whole search space is shown in figure 5.15. The output (f) value 512 yielded the best mean AUC (from the figure 5.15). Though the max AUC obtained by the $f=1024$ and $f=2048$ is also scoring very good. However on this value of f the number of total parameters for the whole network depends as the whole feature map of each DenseNet branch are flattened using flatten then concatenated. Hence multiplied by 2 (since 2 branches). This concatenated feature map is then multiplied by the f size when they gets connected. For example one DenseNet branch has feature map has size= n , concatenated feature map has size= $2n$. Concatenated features gets connected to FC network with output size f results in $2n*f$ parameters increase, and produces output size = f . In other words the smallest f size which gives good performance, is better since it keeps the computations smaller. So $f=512$ is chosen.

5.1.11 Batch size analysis

If the batch size is too low then it takes more time and after a certain size it does not train well too. If the batch size is very big then it may train faster but

they generalize lesser as they tend to converge to sharp minimizers of the training function [22].

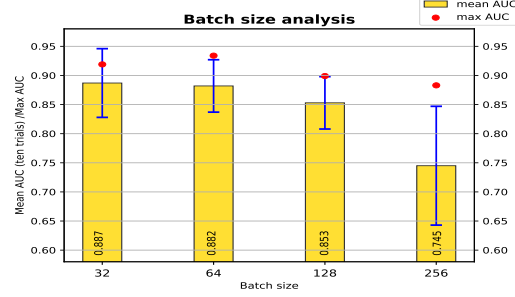


Figure 5.16: DenseNet-Siamese optimal batch size analysis.

Discussion

From the figure 5.16 it is observed that the larger the batch size gets the prediction accuracy on the test data gets worse. Mean prediction AUC for batch size 32 and 64 are very close. The max AUC score for batch size 64 is much higher than 32 and it will train faster too, that is why it is selected as the best value.

5.1.12 Learning rate and optimizer analysis

During the training, in backpropagation step, the analytic gradient is computed which is used to update the parameters of the network (inspired by [34]). This update stage could be done in different ways, this is where the optimizer come into action. While the main target of the deep learning task is to find the minima, the optimizers can control how soon or robustly the minima is found. There is a very compelling comparison of optimization process to a ball or particle rolling down hill in the Stanford lecture series [14]. It compares the loss function to a hill and randomly initializing the network weights to a particle with zero velocity at random points on the hill. Now the optimization process is compared to simulating the particle's motion (parameter vector) of rolling down the hill landscape (loss).

Keras sources [8] gives very brief description of the optimizers. Important optimizers are as follows: **SGD** Stochastic gradient descent optimizer, the very first of it's kind, conceptualized by H. Robbins and S. Munro back in 1951. Even

though it remains one of the most preferred optimizer till date (different variations available e.g., with momentum, Nesterov etc), this optimizer is not evaluated in this work in favor of more theoretically advanced optimizers. In **Adagrad** instead of globally varying the learning rate, the concept of per parameter adaptive learning rate is first introduced by Duchi et al. in Adagrad optimizer. It seems it has a limitation though, the use of monotonic learning rate is often too aggressive and the learning stops too early. This optimizer is also not included in this study in favor of more advanced optimizers. **RMSprop** try to compensate the aggressive monotonically decreasing learning rate from Adagrad by introducing the moving average of squared gradient. **Adam** can be seen as RMSprop with momentum. **Nadam** incorporates Nesterov momentum into Adam. **Adamax** is a variant of Adam which uses infinity norm. **Adadelata** is like Adagrad with moving window of gradient updates.

Optimal learning rate selection is very important for effective learning. However, optimal learning rate varies optimizer to optimizer, hence for learning rate and optimizer a very fine grained search is performed here The search space contains total 20 different learning rates and five optimizers. Each optimizer is evaluated for all 20 learning rates, that makes 100 network configurations which are trained 10 times from scratch for the evaluation and compared on the prediction accuracy (AUC) and std as usual. The search space is presented in table 5.3:

Table 5.3: The search space for learning rate and optimizer best hyperparameters. Since both are related they need to be evaluated together.

Learning rate
0.1, 0.5, 1.0 (only for Adadelata)
0.01, 0.02, 0.03, 0.05, 0.07
0.001, 0.002, 0.003, 0.005, 0.007
0.0001, 0.0002, 0.0003, 0.0005, 0.0007
0.00001, 0.00002, 0.00003, 0.00005, 0.00007
Optimizers
Adam, Nadam, Adamax, RMSprop, Adadelata

The evaluation results are presented in figure 5.17 where the results for different

learning rate are compared for each of the optimizers. In figure 5.18 different representation of the evaluation result is presented which offers the comparative view for each of the optimizers for specific learning rate.

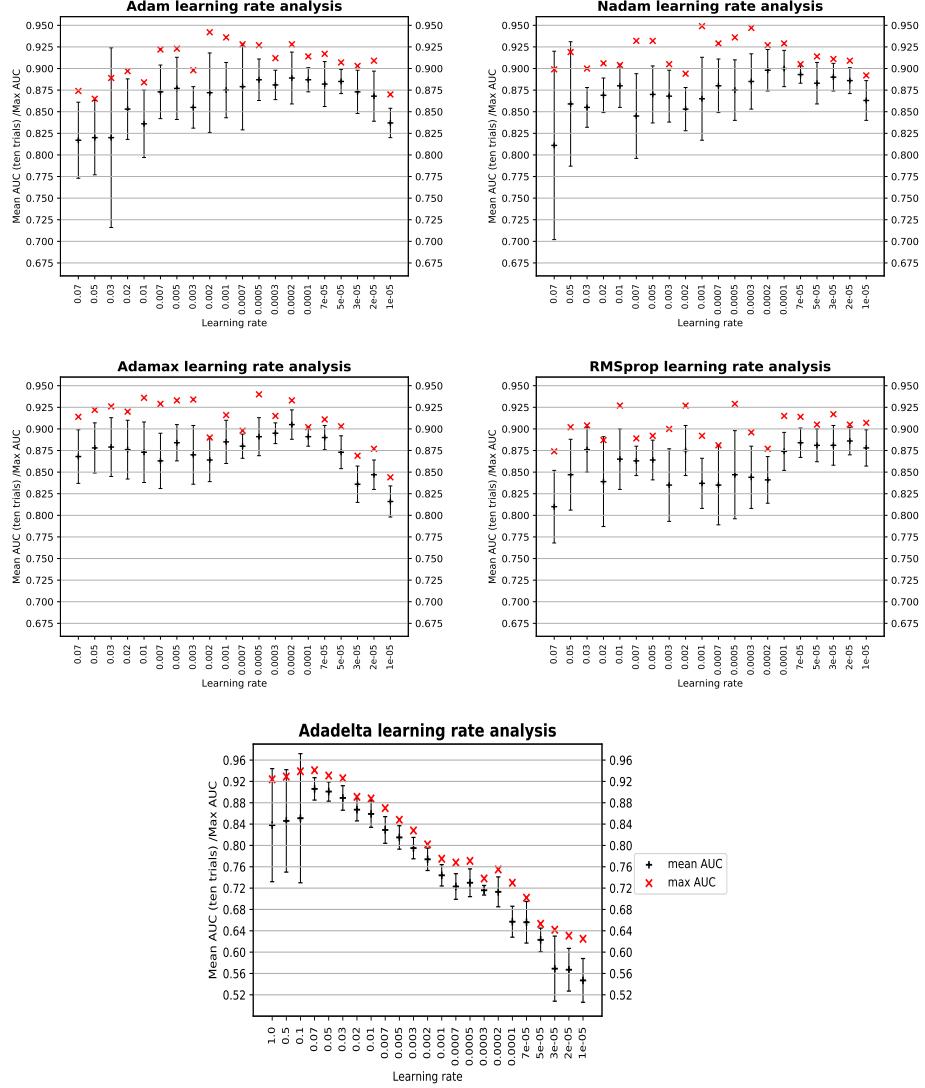


Figure 5.17: In order, (a)Adam, (b)Nadam, (c)Adamax, (d)RMSprop, (e)Adadelata learning rate analysis

Discussion

Overall Adadelata optimizer with learning rate 0.07 has the highest mean AUC

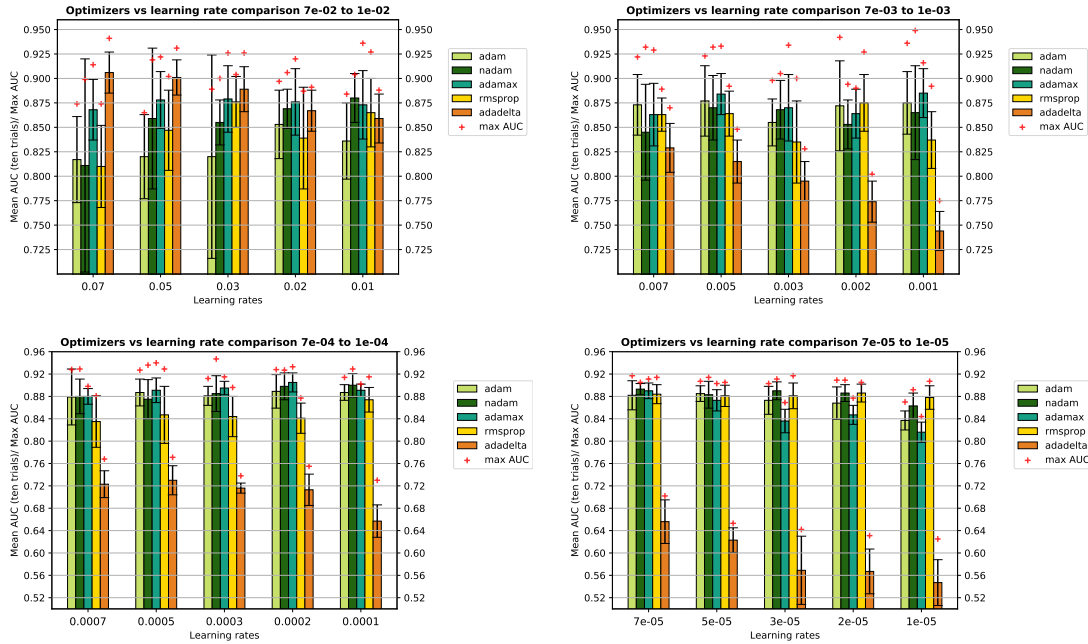


Figure 5.18: Comparison of different optimizers across different learning rates. For more effective comparison the total evaluation is visualized 5 learning rate at a time. Mean AUC obtained for all five optimizers for specific learning rate are grouped together, visualized using bars of different colors. Each two figures in a row have same limits for comparison.

0.906. Which is closely followed by Adamax and Nadam. Adam and RMSprop are slightly behind. The best mean AUC results are displayed in table 5.4. So there is no fixed learning rate for which all the optimizer works best, which is expected as well. For Adamax, Adam, Nadam learning rate 0.0002 works very good, for RMSprop learning rate 2E-5 works better. For Adadelata, as the learning rate drops the performance also drops significantly. For Adadelata the best learning rate is at the boundary condition (0.07) of the evaluation range. So further analysis needs to be done for higher learning rates like 0.1, 0.5, 1.0. In Keras the default learning rate for Adadelata in Keras is 1.0 so it is not really surprising.

Table 5.4: The learning rate for which the optimizers recorded it's best mean AUC.

Optimizer	Learning rate	Mean AUC(10 trials)
Adadelata	0.07	0.906
Adamax	0.0002	0.905
Nadam	0.0001	0.9
Adam	0.0002	0.889
RMSprop	0.00002	0.886

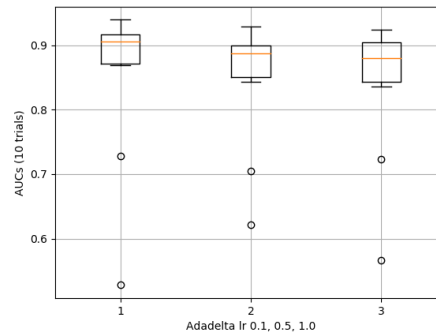


Figure 5.19: Adadelta evaluation for high learning rates. Two of the AUCs out of ten trials for each of the three configurations are displayed as the outlier in the Boxplot. In each of the cases the network did not train well and possibly remained stuck in local minima. The Boxplot displays statistical five number summary for the 10 AUCs recorded for each of the configurations.

From figure 5.19 it is observed that each of the trials with learning rate 0.1, 0.5 and 1.0 somehow has 2 AUCs out of 10, which are far lower than others. Which are displayed in the boxplot as the outliers. After discarding the outliers the mean AUC of the trials would drastically improve and displayed in the following table 5.5. It is however debatable, not fair to calculate mean AUC discarding some runs. In any case best learning rate 0.07 is chosen for Adadelta as it is found to be more dependable and high scoring too.

Table 5.5: Optimizer Adadelta evaluation with high learning rates. Displayed results are after filtering the outlier cases.

Learning rate	0.07	0.1	0.5	1.0
Mean AUC	0.906	0.906	0.891	0.887

5.1.13 Final grid search

The search for the best hyperparameters are done separately so far, now all the best performing hyperparameter values are put together as part of this final grid search. If the hyperparameters that are searched separately, did not have any interdependency (e.g., Learning rate and optimizer has), then the overall prediction accuracy should improve when all the best hyperparameter values are used together.

Hyper-parameters for final grid search

The search space for this final grid search has been narrowed down manifold by doing the individual or focused parameter search before. The remaining search space for the final search is described below in table 5.6.

Table 5.6: Best hyperparameter values for final evaluation.

Name	Value	Name	Value	Name	Value
Number of filter	16	Layers	2-2, 3-4	Growth rate	12, 18, 30
DenseNet dropout	0.2, 0.4	Compression	0.3, 0.7	Bottleneck	'False'
FC output	512	FC dropout	0.5, 0.7	Pooling	'flatten'
Batch size	64	Optimizer & learning rate	'Adadelta' & 0.07, 'Adamax' & 0.0002, 'Nadam' & 0.0001		

From the table 5.6, it is observable that three different optimizers are evaluated with their corresponding learning rates. So basically each test case is evaluated for each optimizer. Similarly for other hyperparameters, which have more than one values in the search space, the overall search cases are multiplied by the number of those many values. Here total of 48×3 (for optimizers) = 144 dimensional search space is being searched for the final evaluation i.e 144 network configurations. Each of them will be trained from scratch 10 times and the network configuration with best mean AUC for test data prediction, is considered to be the best network. And the corresponding hyperparameter values will be the best hyperparameter values for DenseNet-Siamese network.

Table 5.7: Best hyperparameter values for the DenseNet-Siamese, obtained from the final grid search. In decreasing order of their mean prediction accuracy on the test data. The results are very close, hence top 5 results are displayed instead of just one.

Config alias	Epochs	Learning rate	Optimizer	Layers	Growth rate	DenseNet dropout	Compression	Mean AUC	Std	Max AUC
'Config 1'	14	0.07	'Adadelata'	2-2	30	0.4	0.3	0.921	0.016	0.95
'Config 2'	15	0.0002	'Adamax'	2-2	18	0.4	0.7	0.918	0.009	0.935
'Config 3'	14	0.07	'Adadelata'	3-4	12	0.4	0.7	0.915	0.019	0.94
'Config 4'	14	0.07	'Adadelata'	2-2	18	0.4	0.3	0.913	0.012	0.927
'Config 5'	13	0.07	'Adadelata'	2-2	12	0.2	0.7	0.912	0.011	0.932

1097 All the top results were obtained for FC dropout probability 0.7. Other
 1098 hyperparameters which has multiple values in final search space, the best values
 1099 are mentioned in table 5.7. Rest of the hyperparameter values are fixed to what is
 1100 displayed in table 5.6. From table 5.7 the 'Config aliases' are used to refer to the
 1101 network structure in the following figures, for comparing the results.

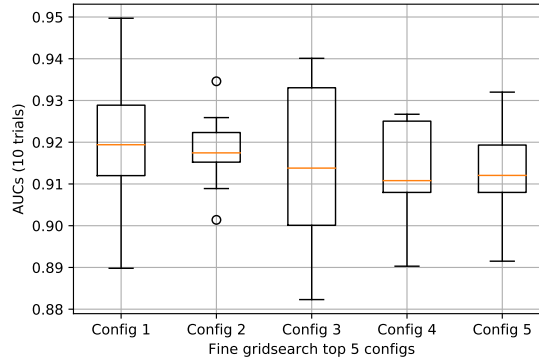


Figure 5.20: Top five DenseNet-Siamese network each evaluated 10 times from scratch. For most effectively compare the network statistical analysis is done for all ten AUC prediction results for each of the network and visualized using Boxplot

1102 Discussion

1103 The best result reported in Valdenegro et al. [41] is **0.91** AUC for binary class
 1104 prediction and **0.894** for the score prediction, for the same data used in this thesis.
 1105 However, In this work only score prediction is done. DenseNet-Siamese network
 1106 able to record highest mean AUC in ten trials **0.921**, with standard deviation
 1107 across the trials of 0.016 and maximum AUC of 0.950. Which is higher than the
 1108 state of the art score prediction. The best performance is recorded with Adadelata



Figure 5.21: Statistical analysis on top results from final grid search visualized in bubbles with color codes which enable the qualitative analysis among different configurations.

optimizer and learning rate 0.07.

The mean AUC result for the top configurations are very close and it is hard to declare a runaway winner as the highest network has mean AUC of 0.921 and std of 0.016, second best has mean AUC of 0.918 with std 0.009.

It is interesting to note that the dropout values which are giving the best results are high 0.7 for the Siamese side and 0.4 (only for 'Config 5' it is 0.2) for the DenseNet side. High dropout is enabling the network to be good in generalization. The training/validation accuracy is could reach very high, still the generalization on test data is high.

The statistical comparison of the evaluations which yielded top 5 mean AUC are displayed in figure 5.20. Each of the 5 network configuration is evaluated for 10 times and the prediction accuracies are compared quantitatively based on the box and whisker plot representation.

Interestingly enough, for 'Config 2', two points are displayed to lie outside the boxes. This supposed outliers, both in higher and lower end of the distribution range, are not really outliers, it just so happens that this two values are bit higher and lower than the rest 8 AUC values, which have only standard deviation of only 0.05. In fact this shows that this result is the most consistent one, it is also the

only top configuration that uses 'Adamax' optimizer and as seen from the 5.7 the reduction ratio is 0.7. In comparison to the 'Config 1' the reduction ratio used for 'Config 2' is higher. As previously seen in figure 5.12 the total parameters for the networks are much lesser for the 'Config 2' than 'Config 1'. So overall this result for 'Config 2' is also very good.

And in figure 5.21 multiple statistical data is displayed in bubbles whose color indicates the scale of the value. Lowest values indicated by brown or shade of red and as it gets higher the color gets orange, yellow, green, cyan to blue and darker blue. Each column may have data in different ranges, so for most effective qualitative comparison each columns, which are one of the statistical measure, the color maps normalized for each of the column. It is done automatically, by assigning the lowest value of the 5 results for each metric to the lowest color map and highest to the highest color map. So in this one image 5.21 results from 50 evaluations are summarized. That is 10 evaluations for each of the 5 network configurations. The visualization displays statistical five number summary (min, Q1, median, Q3, max) and mean and standard deviation are computed from each of the 10 evaluations for each configuration, which signifies how each network performs.

This visualization is intended for rough qualitative analysis among the top configurations. It is clear that the 'Config 1' network outperforms other networks, because overall it has most 'bubbles' which are blue or dark blue, which means high values. Except the min value is very low, it's displayed in brown. Similarly, 'Config 2' is also very good. It has the lowest standard of deviation. Apart from this configuration (with 'Adamax'), other four top configurations are using 'Adadelta'. 'Nadam' has got the best mean AUC of 0.912, which is same as the 'Config 5'. In figure 5.22 the Receiver operating characteristic (ROC) curve is presented for the best network 'Config 1'. The ROC curve is based on the evaluation which has maximum accuracy out of 10 runs for 'Config 1'. With this results the hyperparameter search for DenseNet-Siamese comes to the conclusion.

5.2 DenseNet Two-Channel

In order to get the best result from DenseNet two-channel (**DTC**) network for the dataset used in this work, best hyperparameter values needs to be found out. In the following section how the best hyperparameter search for DTC is conducted

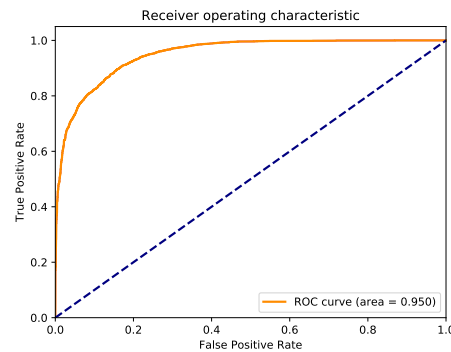


Figure 5.22: The area under curve is visualized for the best network 'Config 1'.

is described. Also the result of the DTC using the best parameter values are evaluated and the results are visualized for easier interpretation.

5.2.1 Hyperparameters to be evaluated

Overall hyperparameters of DTC can be divided into two parts as follows:

Hyperparameters of DenseNet

This important hyperparameters are already mentioned in the section 5.1.1, during the best hyperparameter search for DenseNet-Siamese. While the structure and the intended functional usage are different the DenseNet two-channel and DenseNet-Siamese both depend on the basic feature extraction capability of the network. Hence there are some common knowledge which can be extracted from the hyperparameter search in previous section (DenseNet-Siamese network). But the main hyperparameters such as growth rate, layers per block, number of filters, reduction and bottleneck still needs to be searched again, as the network structure is different.

Other hyperparameters

Apart from the DenseNet there are other general hyperparameters such as learning rate, batch size for training and optimizer.

1178

1179 5.2.2 DenseNet growth rate and layers per block analysis

1180 Similar to the previous analysis, once all the hyperparameters which needs to
 1181 be evaluated are known, manual optimization takes place. After many trial and
 1182 error cycles or randomly the parameter values are changed around to get a intuition
 1183 on the possible ranges for the hyperparameter values to be evaluated. This step
 1184 also ensures that the starting point is not very bad. Layers per block are chosen
 1185 among 2, 3, 4, 5, 6. For single dense block evaluation goes up-to 12 layers (per
 1186 block). Each network has been evaluated for growth rates of 12, 18, 24, 30. Growth
 1187 rate 6, 36 excluded as they are too thin and too thick respectively (found from
 1188 DenseNet-Siamese evaluation). Different dense block sizes of 1, 2, 3, 4. There are
 1189 other parameters but number of dense block, growth rate and layers per block are
 1190 three main parameters. The parameters compression/reduction and bottleneck
 1191 are set 0.5 and False respectively. For more fine grained analysis the reduction
 1192 and bottleneck parameters might be evaluated. `nb_filter` values are fixed at
 1193 16 for this test. The parameter classes are set to 2, where class 1 for matching
 1194 patches and 0 for not-matching patches. 96, 96 is the input image dimensions
 1195 and input is two-channel. So depending on local setting of the keras, channel-first
 1196 or channel-last suitable `input_shape` is chosen automatically as 2, 96, 96 or 96,
 1197 96, 2 respectively. The learning rate used for the test is 0.07 and Adadelta as
 1198 optimizer. Best performing combination from DenseNet-Siamese analysis. Dropout
 1199 for DenseNet used as 0.2 to incorporate minimal regularization. Epochs are different
 1200 for different architectures to ensure that the networks are able to train decently.
 1201 Flatten is used as pooling at the end of the DenseNet, in place of global average
 1202 pooling. Binary_crossentropy loss function with Sigmoid activation function used
 1203 for the binary classification, this final layer acts as the binary classifier. This is
 1204 ensured by including the top of DenseNet architecture (`include_top=True`). In all
 1205 the cases the networks are being trained from scratch. Parameter `weights` value
 1206 None ensures that no previously trained weights are used.

1207 Growth rate and layers per block search setup summary

1208 The overall search space is summarized in the section below.

Fixed hyperparameters

Other hyperparameters that are needed to instantiate the DenseNet are set to fixed values, in order to focus on the layers per block and growth rate parameters.

Table 5.8: Fixed hyperparameter values for the evaluation setup.

Name	Value	Name	Value	Name	Value
Number of filter	16	Subsample initial block	'True'	Weights	'None'
Dropout rate	0.2	Include top	'True'	Compression	0.5
Bottleneck	'False'	Pooling	'flatten'	Transition pooling	'max'
Optimizer	'adadelata'	Learning rate	0.07		

Varying hyperparameters

- **Nb_layers_per_block:**

- **One dense block architecture (nb_dense_block=1)**

'2', '4', '6', '8', '10', '12'

- **Two dense block architecture (nb_dense_block=2)**

'2-2', '4-4', '6-6'

- **Three dense block architecture (nb_dense_block=3)**

'2-2-2', '4-4-4', '6-6-6'

- **Four dense block architecture (nb_dense_block=4)**

'2-2-2-2', '4-4-4-4', '6-6-6-6'

- **Growth rate:**

- Thin layers: 12, 18

- Thick layers: 24, 30

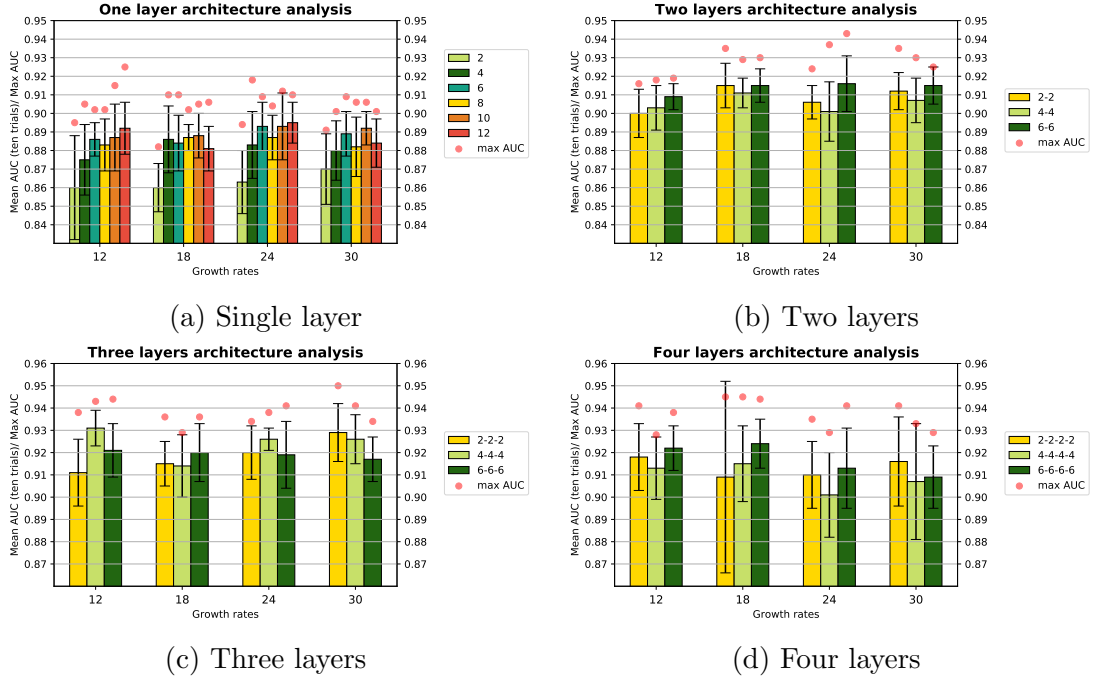


Figure 5.23: DenseNet two-channel layers per block and growth rate analysis. Unique layers per block is represented with a specific color code. Whole result is divided according to dense block numbers, for easier representation. Every two charts presented side by side have identical y-axis for comparison. All the results obtained for same growth rate are grouped together for each of the growth rate (12, 18, 24, 30) represented by x-axis.

Discussion

Each sub-figure shows mean AUC (on the test dataset) and standard deviation of 10 trials and maximum AUC. From figure 5.23 it is clear that 4 dense blocks and 3 dense blocks network performs better than the 1 and 2 block/s networks. The original paper [17] also used 3 and 4 blocks DenseNet for most of the evaluations. So this finding is expected. For growth rates though, no clear trend is observed, in fact, from sub-figure 5.23(c) it is observed that with increase in growth rate the '2-2-2' network performs better, while for '6-6-6' it mildly decreasing and similar for '4-4-4'. So no common best growth rate can not be determined for all the architectures, so it would be safe to evaluate for as many as possible. Growth rate 12, 18 and 30 will be evaluated for the final run.

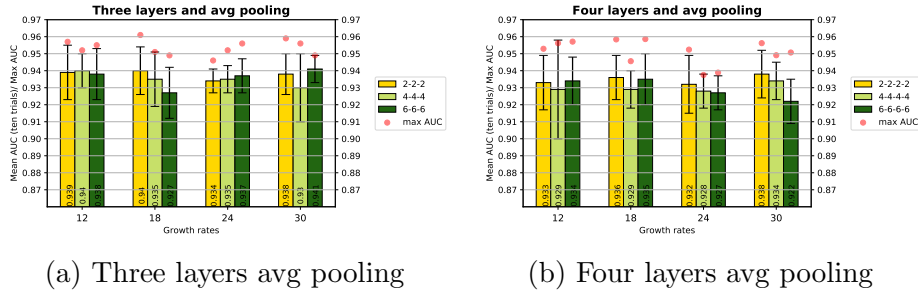


Figure 5.24: DenseNet two-channel avg pooling analysis

5.2.3 Pooling analysis

The type of pooling used at the end of the network also determines the size of the total parameter size and also affects the generalization of the data. The original paper [17] used global average pooling (avg) [27]. In the previous study **flatten** pooling is evaluated, for this analysis **avg** pooling is evaluated for 3 layer and 4 layer DenseNet blocks, since they are the best performing network in previous analysis. In figure 5.24 it is observed that between the three layer and four layers DenseNets, the former performs better. Overall the avg pooling is resulting in smaller total network parameters and also the better mean AUCs.

For a comparative analysis between the flatten and average pooling the mean AUCs are compared for each of the growth rate and for three layers DenseNet (2-2-2, 4-4-4, 6-6-6). Apart this layers and growth rate values, and pooling, other parameters remain same as the basic evaluation network 5.2.2.

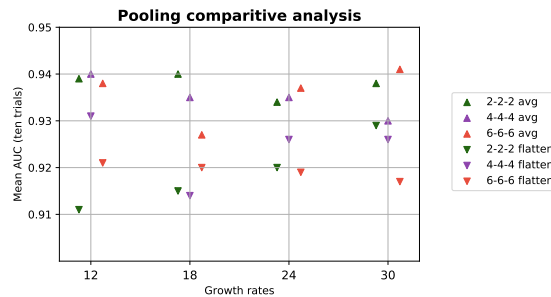


Figure 5.25: Average vs flatten pooling

Discussion

Here, the triangle up represents the mean AUC obtained using average pooling, triangle down represents flatten pooling. The results from the DenseNet architectures have been grouped together for each growth rates accordingly. So the x-axis of the graph shows growth rates. The y-axis represents the mean AUC obtained from 10 trials of each configurations. From figure 5.25 it is clearly seen that the average pooling produces better results than using flatten, in all the cases. Since the representation might be bit complex, the same data from figure 5.25 is displayed in tabular view (5.9) as follows.

Table 5.9: The average and flatten pooling comparison results.

Architecture	Growth rate	Mean AUC (Flatten)	Mean AUC (Average)
2-2-2	12	0.911	0.939
2-2-2	18	0.915	0.94
2-2-2	24	0.92	0.934
2-2-2	30	0.929	0.938
4-4-4	12	0.931	0.94
4-4-4	18	0.914	0.935
4-4-4	24	0.926	0.935
4-4-4	30	0.926	0.93
6-6-6	12	0.921	0.938
6-6-6	18	0.92	0.927
6-6-6	24	0.919	0.937
6-6-6	30	0.917	0.941

5.2.4 Basic network structure

For doing effective hyperparameter search, need to evaluate the target parameter with the predefined set of values, while all the other parameters which are not directly related are kept constant. So for this purpose a common network structure is defined, which is name 'basic network'. For all the hyperparameter searches, this structure will be the same and only the target hyperparameter values will change.

Table 5.10: DenseNet two-channel 'basic network structure' for further hyperparameter search.

Architecture	Growth rate	Nb_filter	Dropout	LR.
2-2-2	30	16	0.2	0.07
Optimizer	Reduction	Bottleneck	Batchsize	Pooling
adadelata	0.5	FALSE	64	avg

1266

1267

5.2.5 Number of filter analysis

1268 Initial number of filters. 8, 16, 32, 64 are being evaluated here. keeping all
 1269 other parameters fixed, a comparison between mean AUCs (10 trials) obtained with
 1270 growth rate 18 and 30 is also done under this analysis. To explore if the growth
 1271 rate and number of filter have any effect on each other.

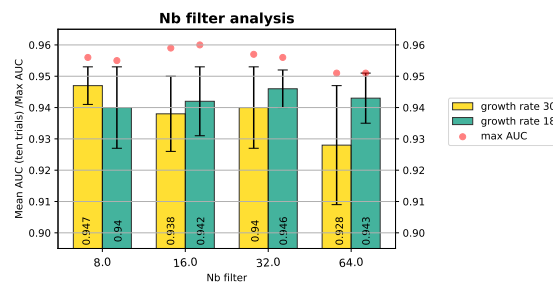


Figure 5.26: DTC Number of filter size analysis

1272

Discussion

1273 It is observed in 5.26 that with change in number of filter size the mean AUC
 1274 varies a lot for higher growth rate such as 30, for growth rate 18 it does not vary
 1275 so much. For growth rate 30 the nb filter 8 has the best mean AUC. For growth
 1276 rate 18 the nb filter 32 has the best mean AUC. Overall growth rate 30 with nb
 1277 filter 8 and growth rate 18 with nb filter 32 are the best combinations.

1278

1279 **5.2.6 Reduction and bottleneck analysis**

1280 This analysis is for evaluating the effect of different reduction ratios and the
 1281 effect of bottleneck. So the mean AUC is recorded for 10 trials for each of the
 1282 reduction values 0, 0.2, 0.3, 0.5, 0.7. This is a rather coarser search space. But
 1283 each of them is also evaluated with bottleneck as well, the effect of varying values
 1284 of reduction and with/without bottleneck is displayed in the figure 5.27.

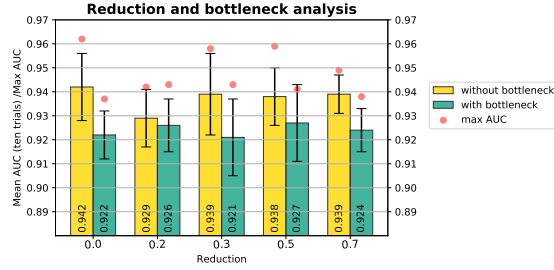


Figure 5.27: Reduction and bottleneck analysis

1285 **Discussion**

1286 The effect of the bottleneck layer is rather limiting the generalization of the
 1287 network. So it is observed that without bottleneck should be used for future
 1288 evaluations. The performance without reduction is best as expected, how ever main
 1289 purpose of the reduction is to decrease the number of total parameters. The mean
 1290 AUC are highest and almost same for 0.3, 0.5 and 0.7. The original implementation
 1291 [17] uses 0.5 as reduction ratio. Here also 0.5 is selected as the best reduction ratio
 1292 for the final evaluation.

1293

1294 **5.2.7 Total parameters analysis**

1295 For the total parameter comparison 2, 2-2, 2-2-2, 2-2-2-2 layers per block
 1296 networks are chosen along with growth rate 18. Every other configurations are
 1297 kept same, with avg pooling and without any reduction. With same configurations,
 1298 except growth rate, the relation between growth rate and total parameters are
 1299 evaluated.

Table 5.11: Total number of parameters are compared across different number of dense block and different growth rates.

Layers per block	Growth rate	Total parameters
2	18	17217
2-2	18	47857
2-2-2	18	96785
2-2-2-2	18	166593
2-2-2	12	55529
2-2-2	18	96785
2-2-2	24	149201
2-2-2	30	212777

From 5.11 it is noted that the total parameters increase as the number of dense blocks increase. Same trend is noted for growth rates too. Total parameters also varies with the choice of pooling and reduction, a simple example of each cases is as follows. Keeping all other configuration same, network has 96785 parameters when `avg` pooling is used, and 101685 is the size with `flatten` pooling, which is more parameters. Keeping all other configuration same, with 0.5 reduction ratio for DenseNet the network size goes down to 51430 instead of 96785 total parameters without reduction. Both are evaluated for 2-2-2 layers and growth rate 18.

5.2.8 Dropout probability analysis

10 evaluations are done for each test cases, consisting of values ranging between 0 and 0.7, with an interval size of 0.1.

Discussion

In figure 5.28 it is observed that the 0.2 dropout configuration obtained the highest mean AUC. The other values with lesser dropout probability or greater dropouts are all gradually decreasing as they go further from the peak (0.2). With exception of the mean AUC obtained with 0.7 dropouts. The best value selected for further evaluation is 0.2.

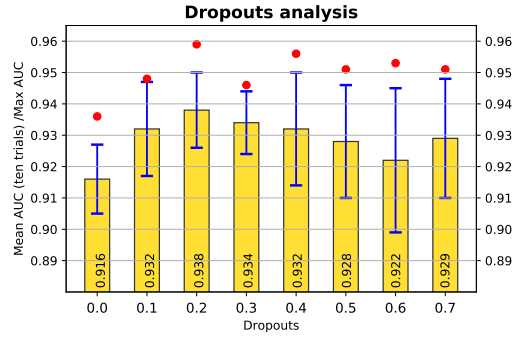


Figure 5.28: DTC Dropout probability analysis

1318

1319 **5.2.9 Batch size analysis**

1320 For this analysis a limited search space is defined, with only batch size varies,
 1321 the evaluation (10 trials for each test case) results are displayed in figure 5.29.

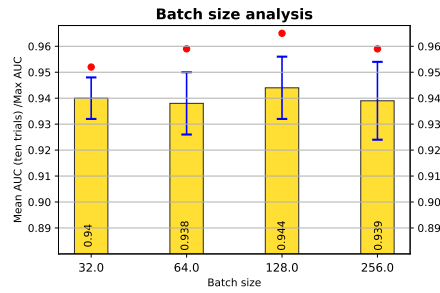


Figure 5.29: Batch size analysis

1322 **Discussion**

1323 From the figure 5.29 it is concluded that the batch size of 128 works the best.

1324

1325 **5.2.10 Learning rate and optimizer analysis**

1326 For this analysis only Adadelta optimizer is used. This is based on the intuition
 1327 that is formed during the DenseNet-Siamese evaluation. But for the optimal
 1328 learning rate the a series of discrete values are evaluated. The search space and
 1329 results are presented in figure 5.30.

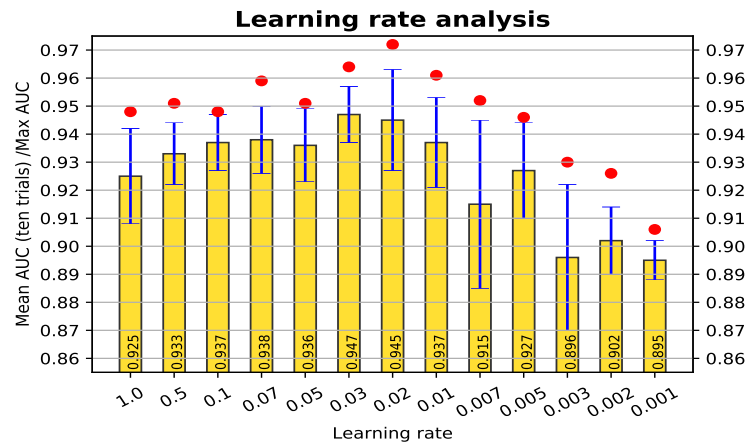


Figure 5.30: Learning rate analysis

Discussion

From figure 5.30 it is observed that the mean AUC with the learning rate 0.03 is slightly higher than the others. While one of the evaluation with 0.02 learning rate has obtained the maximum AUC of 0.972. But 0.03 is chosen as the best learning rate for further evaluations.

5.2.11 DTC final grid search

Table 5.12: Best hyperparameter values for final evaluation of DenseNet two-channel.

Name	Value	Name	Value
Layers	2-2-2, 4-4-4, 6-6-6	Pooling	'avg'
Growth rate (gr) & Number of filter	12 & 32, 18 & 32, 30 & 8		
DenseNet dropout	0.2	Compression	0, 0.5
Bottleneck	'False'	Batch size	128
Optimizer & learning rate	'Adadelta' & 0.03		

Layers per block 2-2-2 it has very consistent performance in terms of mean AUC and also able to score high max AUC. It could have been possible to do architecture searches for 3 dense block architectures with layers per block close to 2-2-2, for example 2-3-3 etc. But then the search grid will be very big. Also

Table 5.13: Final grid search results

Gr.	Metrics	Layers per block					
		2-2-2		4-4-4		6-6-6	
		R=0	R=0.5	R=0	R=0.5	R=0	R=0.5
12	Mean AUC	0.95	0.944	0.95	0.95	0.947	0.945
	Std	0.011	0.015	0.009	0.01	0.008	0.008
	Max AUC	0.97	0.97	0.963	0.965	0.963	0.955
	Total Parameters	55,529	30,163	159,473	87,629	317,561	176,535
18	Mean AUC	0.952	0.955	0.951	0.944	0.948	0.938
	Std	0.008	0.009	0.005	0.011	0.006	0.014
	Max AUC	0.967	0.966	0.956	0.963	0.956	0.955
	Total Parameters	96,785	51,430	308,369	168,671	640,481	355,860
30	Mean AUC	0.943	0.948	0.943	0.944	0.932	0.941
	Std	0.008	0.008	0.01	0.013	0.015	0.011
	Max AUC	0.959	0.964	0.96	0.962	0.948	0.953
	Total Parameters	160,001	82,162	650,873	355,949	1,473,665	822,276

included in the evaluation, layers 4-4-4 and 6-6-6. All other hyperparameters are displayed in the table 5.12

The result is displayed in the table 5.13. unlike Densenet-Siamese this final search space is much smaller and whole search space is visualized here in tabular format. Three architectures 2-2-2, 4-4-4 and 6-6-6 are evaluated for all three growth rates 12, 18, 30 and also for Reduction 0.5 and without Reduction. In the table 5.13 the growth rate is displayed as Gr. and Reduction is displayed as R for space constraint.

Discussion

The best result obtained has mean AUC of **0.955**. This is with reduction 0.5, 2-2-2 layers per block and growth rate of 18. Along with the other values mentioned in 5.12 these are the best hyperparameters for DenseNet two-channel, which might be specific to the dataset for the thesis. Normally it is observed that the 2-2-2 performance is very similar to that of 4-4-4, in fact slightly better. The performance of 6-6-6 is bit worse than the other too. Although, because of

reduction the auc is observed to be slightly lower some times, some times it is higher than the without reduction result. But the size of the total parameters of the network with $\text{Reduction}(R)=0.5$ is always close to half size of the equivalent network without Reduction. So that is always beneficial as it is less computationally expensive. In Valdenegro et al. [41] work it is also found that the simple two-channel network better than the Siamese network. By comparing DenseNet two-channel and DenseNet-Siamese same trend is observed.

During the hyperparameter search many high accuracy network models are recorded, the highest AUC recorded for any of the DenseNet two-channel is 0.972, which is found with 'basic network' structure and learning rate 0.02. The ROC for the prediction is displayed in figure 5.31.

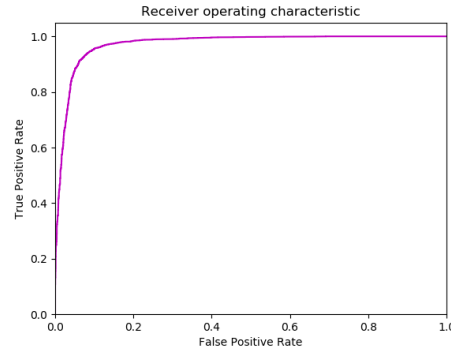


Figure 5.31: Overall best result for DenseNet two-channel 0.972 AUC (Single run)

Top 6 configurations

Instead of top 5, top 6 configurations are mentioned here, because the 6th one has same mean AUC and 4th and 5th. Table 5.12 contains the overall search space for this search. Those parameters which has multiple values in the search space, best values for those are mentioned table 5.14, e.g. reduction, 0 and 0.5. In table 5.14 the 'Config alias' are the aliases for the configurations, for future reference.

Discussion

From the boxplot 5.32, how the 10 prediction values on test data are distributed are evaluated, with a statistical outlook. For example the distribution for third and

Table 5.14: Best hyperparameter values for the DenseNet two-channel, obtained from the final grid search. In decreasing order of their mean prediction accuracy on the test data. Top 6 results are mentioned.

Config alias	Layers	Growth rate	Epochs	Compression	Mean AUC	Std	Max AUC
'Config 1'	2-2-2	18	'21'	0.5	0.955	0.009	0.966
'Config 2'	2-2-2	18	'25'	0	0.952	0.008	0.967
'Config 3'	4-4-4	18	'21'	0	0.951	0.005	0.956
'Config 4'	2-2-2	12	'23'	0	0.95	0.011	0.97
'Config 5'	4-4-4	12	'21'	0	0.95	0.009	0.963
'Config 6'	4-4-4	12	'21'	0.5	0.95	0.01	0.965

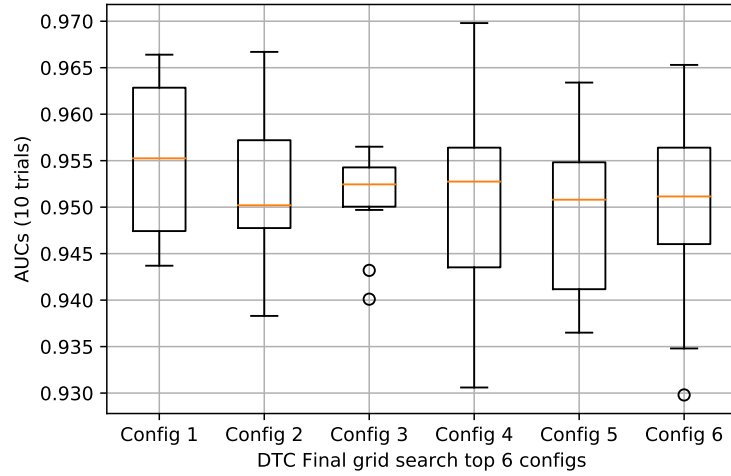


Figure 5.32: Boxplot analysis for top configurations, with statistical outlook on 10 AUC scores obtained for each cases.

second configuration are very condensed, i.e less standard deviation, which is also supported from the data in table 5.14. Two of the configurations display outliers, which means those readings are much wider in comparison to other readings. But the 'Config 2', which has overall second highest average AUC, has the lowest median. This signifies most data points are actually lower than the mean of the distribution, just few high auc scores pulled up the overall average. So with respect to median it is the worst configuration among the 6. From table 5.14 it is observed that growth rate 18 is the best performer, followed by 12, while only two out of six entries are using compression. 2-2-2 is best performing network, but 4-4-4 is also very close. Even the median for 'Config 1' is highest, along with highest mean AUC, this is the overall best DenseNet two-channel configuration.

5.3 VGG-Siamese network with Contrastive loss

In this section, important hyperparameters of VGG-Siamese network with contrastive loss are listed and the grid search result is presented with help of visualizations. At the end, the optimized prediction result with best hyperparameter values obtained from the grid search is also presented.

5.3.1 Hyperparameters to be evaluated

The hyperparameters can be divided into two types, VGG branch related and other general hyperparameters. In Siamese part the distance is computed directly from the output of the branches. As decision, the network outputs the computed distance between the input patches, so the decision network of the Siamese does not have any hyperparameters. Each VGG branch has the architecture as follows: Conv(n , $a \times a$)-Conv(n , $a \times a$)-MP(2, 2)-Conv($2n$, $a \times a$)-Conv($2n$, $a \times a$)-MP(2, 2)-Conv($4n$, $a \times a$)-Conv($4n$, $a \times a$)-Conv($4n$, $a \times a$)-MP(2, 2)-Conv($8n$, $a \times a$)-Conv($8n$, $a \times a$)-Conv($8n$, $a \times a$)-MP(2, 2)-Conv($8n$, $a \times a$)-Conv($8n$, $a \times a$)-Conv($8n$, $a \times a$)-MP(2, 2)-FC(d). There could be up to 3 FC layers at the end, denoted by l . Also there could be batch normalization layers after all the FC layers in VGG branch. Variables n , a , l , d are the hyperparameters of the VGG branches. The overall search space for the hyperparameters are displayed in the table 5.15.

Table 5.15: The search space for VGG-Siamese network with Contrastive loss.

Name	Value	Name	Value
Conv filters	8, 16, 32, 64	Kernel size	3, 5, 7
FC Layers	1, 2, 3	FC output	32, 64, 96, 256, 512, 1024, 2048
Batch normalization	True, False	Dropout	0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7
Batch size	32, 64, 128, 256, 512	Optimizer	Adam, Nadam, Adadelata, Adamax, RMSprop
Initializers	He_normal, He_uniform, Glorot_uniform, Glorot_normal, RandomNormal		
Learning rates	0.01, 0.007, 0.005, 0.002, 0.0007, 0.0005, 0.0002, 0.0001, 0.00007, 0.00005, 0.00002, 0.00001		

Overall search space is divided into smaller parts and evaluated with a common basic network structure, that gives good prediction accuracy already. The Basic network structure is obtained after manual tuning and many trials.

Table 5.16: Basic network configuration for the hyperparameter search. In this section of report this configuration will be referred to by the alias of 'basic network structure'.

Conv filters	Kernel size	FC Layers	FC output	Batch normalization
16	3	1	128	FALSE
Dropout	Batch size	Optimizer	Initializers	Learning rates
0.2	64	adam	random normal(Conv) + He normal(FC)	0.00001

5.3.2 Basic network structure

This network structure is used for all the focused hyperparameter search. Keeping other values same, only changing the values for the targeted and related parameters, should help building an intuition of how the network behaves with different values of the targeted parameter and should help in selecting the best value from the predefined search space.

5.3.3 Flipped labels

Since contrastive loss returns projected distance, here close to zero means similarity and 1 means dissimilarity. Although, in our original data label 1 represents similarity between patches. Hence the labels for train, validation and test data here are all flipped. Operation `new_label = 1 - old_label` is applied to all three data labels. So for this evaluation input label 0 means similarity between patches.

5.3.4 Conv filters analysis

The `filters` [5] defines the number of output filters in each convolution layers. Now for all the 13 convolution layers in the network the filters size can be easily calculated from the first filter size. In previous section it has been shown that the filters for the Conv layers are n , n , $2n$, $2n$, $4n$, $4n$, $4n$, $8n$, $8n$, $8n$, $8n$, $8n$, $8n$ respectively. Here $4n$ means 4 times n . So for each value of filters (n) in the search space [8, 16, 32, 64], the filter values are computed for all 13 Conv layers and each configuration is evaluated 10 times. Evaluation results are displayed in the figure 5.33.

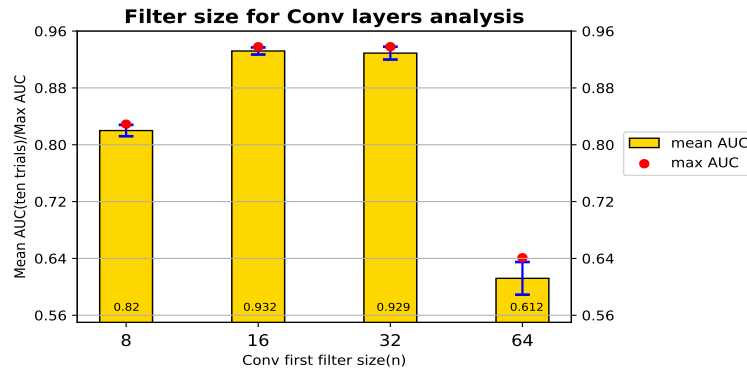


Figure 5.33: Conv filters size analysis for VGG branches.

Discussion

From figure 5.33 the highest mean AUC is for 16. So for the best performing network is having the VGG branch with filters for the convolution layers as follows, 16, 16, 32, 32, 64, 64, 64, 128, 128, 128, 128, 128, 128. The performance for filter size 32 is also very good, but 16 is selected as the best value.

5.3.5 Kernel size analysis

The kernel size parameter defines the height and width of the 2 dimensional convolution window for all the Conv layers in the VGG network. From table 5.3.5 it is clear that with kernel size 3 network learns much better than 5 and 7 kernel sizes.

Rank	Kernel size	Mean AUC (10 trials)
1	3	0.932
2	5	0.799
3	7	0.481

5.3.6 FC units size analysis

The hyperparameter FC output, also known as units in Keras, denoted by (d), determines the output size of the layer. For each values of (d) from the search space of [64, 96, 128, 256, 512, 1024, 2048], one, two and three FC(d) layers along

with ReLU activations are added at the end of each VGG branch. Ten trials are done for each of the configuration. So the goal of this grid search is to find the best values of d and l . The search results are displayed in the figure 5.34.

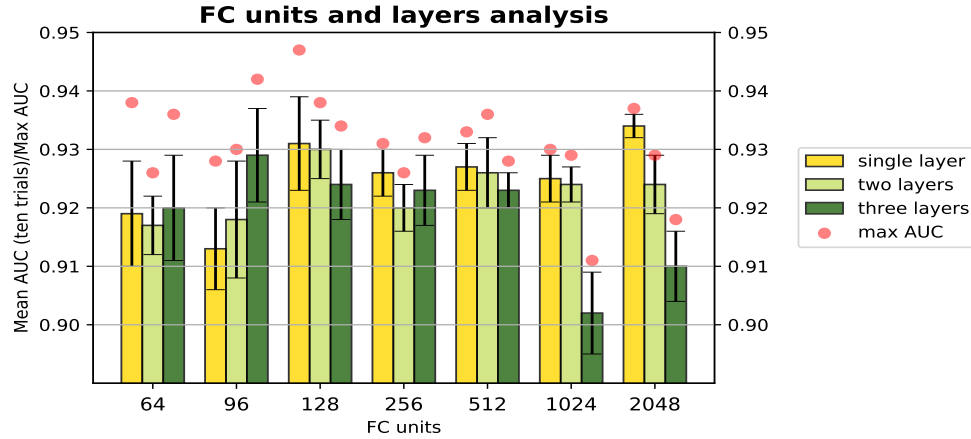


Figure 5.34: FC units and layers analysis for VGG branches.

Discussion

From 5.34 single FC layer ($l=1$) in figure works best. Although performance of two FC layers ($l=2$) are also very close, single layers are selected as best values. Just for validation, the best configuration of the final grid search can also be verified for two layer structures. 128 and 2048 has been selected for the final grid search as they have the highest mean AUCs. Although this is bit unexpected to see two points close to the opposite extremes of search space to yield best results.

5.3.7_INITIALIZER

Keras initializers[7] control or define the way the initial random weights in keras layers are set. In the following part the brief introduction of the initializers and how they can be instantiated in Keras is presented.

1. **Zeros** : This is one of the simplest of the initializers. It generates tensors initialized to 0. In Tensorflow, tensors are represented by n-dimensional arrays of base data type.

Instantiation: `keras.initializers.Zeros()`

1471 2. **RandomNormal** : This initializer uses normal distribution to initialize
1472 generated tensors.

1473 Instantiation: `keras.initializers.RandomNormal(mean=0.0, stddev=`
1474 `0.05,seed=None)`

1475 3. **RandomUniform** : This initializer uses uniform distribution to initialize
1476 generated tensors.

1477 Instantiation: `keras.initializers.RandomUniform(minval=-0.05,`
1478 `maxval=0.05, seed=None)`

1479 4. **Glorot_normal** : Glorot normal initializer is also known as Xavier normal
1480 initializer. It generates samples from a truncated normal distribution which
1481 is centered at 0 with standard deviation of `(stddev) = sqrt(2 / (fan_in`
1482 `+ fan_out))`. `fan_in` represents number of input units in the weight tensor.
1483 The number of output units in the weight tensor is denoted by `fan_out`.

1484 Instantiation: `keras.initializers.glorot_normal(seed=None)`

1485 5. **Glorot_uniform** : Glorot uniform initializer is also called Xavier uniform
1486 initializer. It draws samples from a uniform distribution within `[-limit, limit]`
1487 where `limit` is `sqrt(6 / (fan_in + fan_out))`, `fan_in` represents number
1488 of input units in the weight tensor. The number of output units in the weight
1489 tensor is the `fan_out`.

1490 Instantiation: `keras.initializers.glorot_uniform(seed=None)`

1491 6. **He_normal** : It draws samples from a truncated normal distribution centered
1492 on 0 with standard deviation `(stddev) = sqrt(2 / fan_in)` where `fan_in`
1493 represents the number of input units in the weight tensor.

1494 Instantiation: `keras.initializers.he_normal(seed=None)`

1495 7. **He_uniform** : He uniform variance scaling initializer draws samples from a
1496 uniform distribution within `[-limit, limit]`. Here, `limit` is `sqrt(6 / fan_in)`
1497 and `fan_in` represents the number of input units to the weight tensor.

1498 Instantiation: `keras.initializers.he_uniform(seed=None)`

1499 Kernel initializers for the convolution layers and kernel initializers for the FC
1500 layers are investigated here. The search space for kernel initializers for convolution

layers (Conv) contains [He normal and uniform, Glorot normal and uniform and random normal]. These initializers are selected after some manual trials. With random uniform (default Instantiation) as Conv initializer, the network fails to converge. Hence it is not evaluated. The search space for initializers for FC layer contains [He normal and uniform and Glorot normal and uniform]. The bias initializer, for both Conv and FC layers, is left with the default Zeros initializer. In keras, for both Conv and FC the default kernel initializer is Glorot uniform. Popular intuition is that Glorot or Xavier initialization works better with Sigmoid activation, while He uniform/normal works better with ReLU. In the network architecture there are 13 Conv layers compared to only one or two FC layer/s, so the Conv initializer is expected to have more effect. All the initializers used with default instantiation.

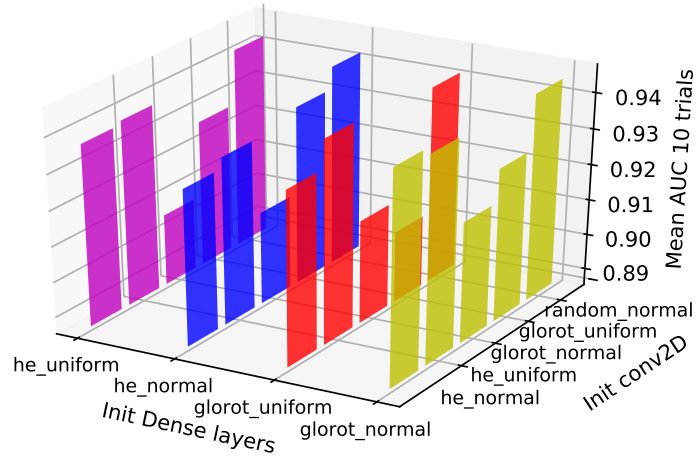


Figure 5.35: Performance comparison of combinations for Conv and FC layer initializers, in one axis the initializers for FC or Dense layers are displayed in unique colors, in the other axis the initializers for Conv layers are showed. Z-axis or height of the bars represents the mean AUC in ten trials obtained for each combination of Conv and FC initializers.

Rank	Init Conv	Init FC	Mean AUC	STD	MAX AUC
1	He normal	Glorot normal	0.943	0.007	0.953
2	He uniform	Glorot normal	0.941	0.012	0.961
3	Random normal	Glorot normal	0.941	0.004	0.946
4	He uniform	Glorot uniform	0.94	0.008	0.952
5	Random normal	He normal	0.939	0.005	0.946
6	Random normal	He uniform	0.939	0.005	0.948

Table 5.17: Kernel initializer top results

Discussion

In this figure 5.35, observed mean AUC is reported for different combinations of Conv initializer and FC Initializer. The x-axis of the graph represents FC layer initializer and the y axis represents the Conv initializer. In the z axis the mean AUC of ten trials is shown. The z-axis values are clipped and starts from 0.89. This is to be able to show the differences in values (bar heights) effectively. Otherwise, if the bars are plotted from zeros the difference is hard to perceive as the values are very close. Lowest mean AUC value obtained is 0.903 and the highest is 0.943. The bar chart starting value selected in such a way that all the bars are visible and comparable.

The performance of random normal as the initializer for Conv layers is over all good. Performance of the Glorot normal and uniform as Conv initializer is comparatively worse than others. Performance of both He initializers are over all good as Conv initializer. As the FC layer initializer glorot normal is found to have performed the best. Over all performance wise the combination of He normal as the Conv and glorot normal as the FC initializer, have yielded highest mean AUC (0.943 AUC). The results for top 6 combinations are presented in table 5.17 in descending order of the yielded mean AUC.

A strong trend is observed with random normal Conv initializers, the standard deviation for it is very low, though the maximum AUC values are consistently lower than the He variations. Glorot normal works very consistently as initializer for FC layers. Therefore it is selected as the FC initializer. But there is no clear

winner for Conv initializers, so both He and random normal initializers are selected for the final grid search.

5.3.8 Batch normalization analysis

Small batch training is better than one by one training and it is also better than training the whole dataset at once. Small batch training with batch normalization is advantageous because it converges faster than doing one by one. Batch normalization reduces the need for carefully tuning the initial weights, also to some extent limits the need for too much of regularizers such as dropouts. Concept of batch normalization is introduced by Ioffe and Szegedy in 2015 [21]. The authors are influenced by the idea from Lecun,1998b [25] and Wiesler and Ney, 2011 [43] that if inputs to a CNN are linearly transformed to have unit variance and zero mean, then the network will converge faster. The learning rates are generally kept comparatively lower because an outlier might cause big effect in already learned activations. As a result of keeping the inputs normalized the outlier cases also affects the overall learning process lesser. Hence batch normalization should also enable the use of higher learning rates. The batch normalization layer is usually applied after each fully-connected or dense layer, apart from the output. In this thesis, the batch normalization is evaluated three steps, firstly, without batch normalization. Then, with batch normalization, when adding the batch normalization layer after the dense layer but before the activation function 'ReLU'. Thirdly, after both, dense layer and activation 'ReLU'. Mean AUC comparisons are as follows (same epochs).

Table 5.18: Batch normalization analysis for VGG branches.

Rank	Batch normalization	Mean AUC (10 trials)	STD	Max AUC
1	Without	0.932	0.005	0.938
2	Before activation ReLU	0.9	0.012	0.926
3	After activation ReLU	0.874	0.011	0.893

The results from 5.18 are bit surprising. It is noted that with batch normalization during the training the network can achieve higher training accuracy in the same

number of epoch than without batch normalization. However the generalization performance on the test data is worse using batch normalization. With batch normalization cases are also tested for lesser and greater epochs than without batch normalization, it still scores poorer results.

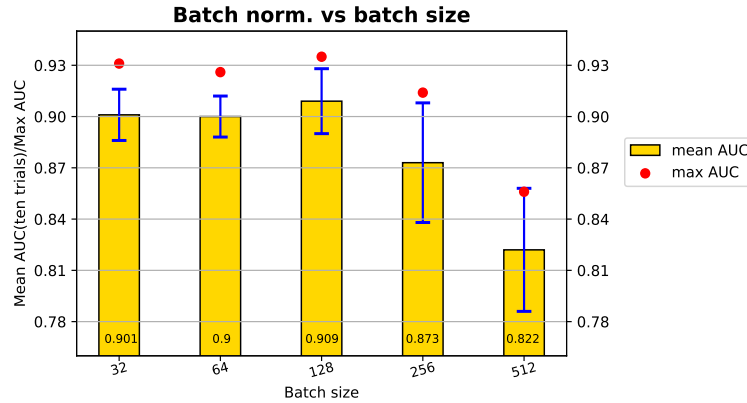


Figure 5.36: Evaluation of how batch normalization affects prediction with increase in batch size.

The batch size and batch normalization correlation analysis is done in figure 5.36. The idea is that the batch normalization performance might increase with increase in the batch size. The bigger the batch size is, the more it resemble the actual distribution that represents the whole data, this is one of the intuition. No strong evidence are found in this direction, though the mean AUC slightly peaked from batch size 32 to 128 and then sharply fell down.

Since batch normalization limits the effect of outliers, could it be the case that the apparent outliers contained discriminative features somehow. Because of sonar images have low signal to noise ratio, it could be that the perceived noise is actually useful data, which get somewhat filtered out by the batch normalization. However this is only speculation, no concrete analysis is done in this direction. But batch normalization layers are not used for further evaluation.

5.3.9 Dropout probability analysis

From the figure 5.37 it is noticeable that the performance (mean AUC across 10 trials) for dropout values 0.3, 0.4, 0.5 and 0.6 are very close. The network achieved

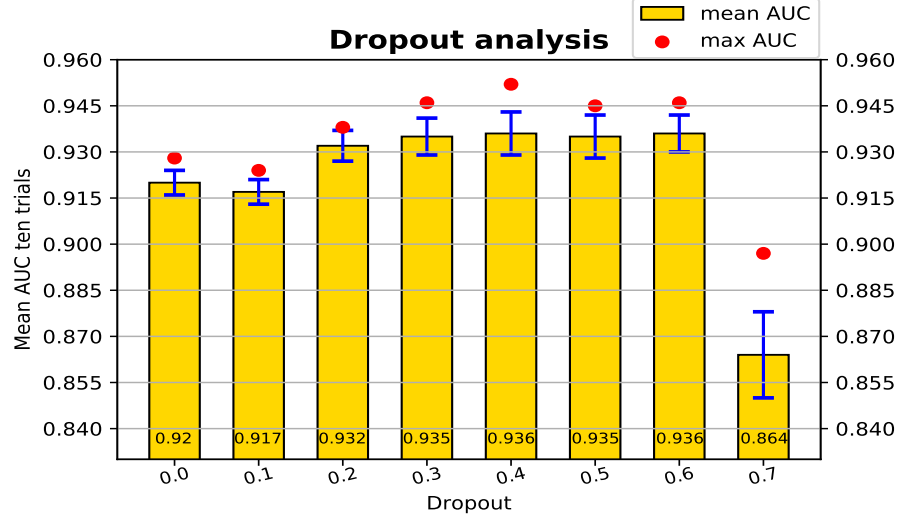


Figure 5.37: Dropouts probability analysis for VGG branches.

1581 highest mean AUC with dropout probability 0.4 and 0.6, so they are chosen as the
 1582 best values. With those values network performed better than dropout 0 and 0.1,
 1583 which offers too less regularization leading to overfit. Also performed better than
 1584 0.7 dropout which provides too much regularization. This is a expected result, and
 1585 matches on the common experience.

1586

1587 5.3.10 Learning rate and optimizer

1588 The search space for optimizer contains Adam, Nadam, Adamax, RMSprop,
 1589 Adadelta. For the learning rate, evaluation started with learning rate search space
 1590 as [0.007, 0.005, 0.002, 0.001, 0.0007, 0.0005, 0.0002, 0.0001, 0.00007, 0.00005,
 1591 0.00002, 0.00001] for all the optimizers. It is noticed that the network, for all
 1592 optimizers, does not train for all the learning rates. Finally, learning rates for each
 1593 optimizer is selected individually. The search space is displayed below in table 5.19.

1594 Discussion

1595 Ten trials per test cases are done. The learning rate for which each optimizer
 1596 obtained best mean AUC result, are compared and visualized in figure 5.38, it is
 1597 noted that the best mean AUC for all the optimizers are very good, though for
 1598 different learning rates. For each optimizer the best learning rates are as follows:

Table 5.19: The search space for learning rates specific for each optimizer for VGG-Siamese and Contrastive loss network.

Adamax	Adadelata
0.001, 0.002, 0.005, 0.007	0.1, 0.2, 0.5, 0.7, 1.0
0.0001, 0.0002, 0.0005, 0.0007	0.01, 0.02, 0.05, 0.07
0.00001, 0.00002, 0.00005, 0.00007	0.005, 0.007
Adam, Nadam and RMSprop	
0.0001, 0.0002, 0.0005	
0.00001, 0.00002, 0.00005, 0.00007	

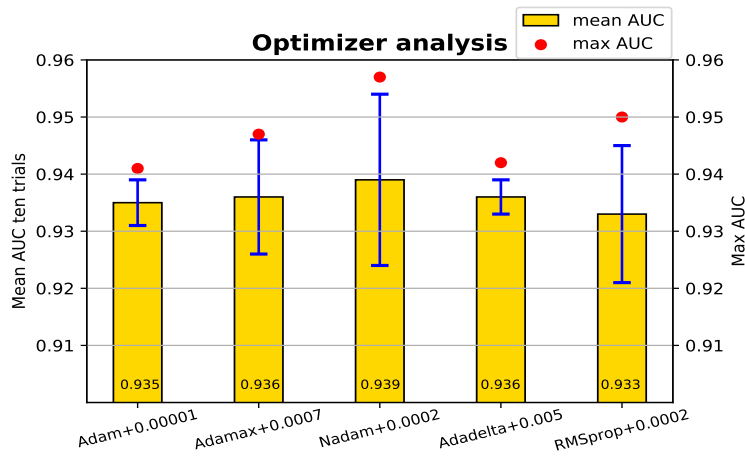


Figure 5.38: Search results of best optimizer and learning rate

Adam-0.00001, Nadam-0.0002, Adadelata-0.005, Adamax-0.0007, RMSprop-0.0002. Nadam is the best performer. From figure 5.38 network with Nadam has highest mean AUC and highest max AUC both. So Nadam optimizer and starting learning rate 0.0002 is selected for the final grid search.

5.3.11 Batch size analysis

To find the best optimal batch size, a search space containing [32, 64, 128 and 256] is evaluated.

It is observed from figure 5.39 that the test AUC somehow increases with the increase in batch size, and opposite to expectation, takes more epochs to reach the convergence. Since the 256 batch size is at the boundary of search space, which

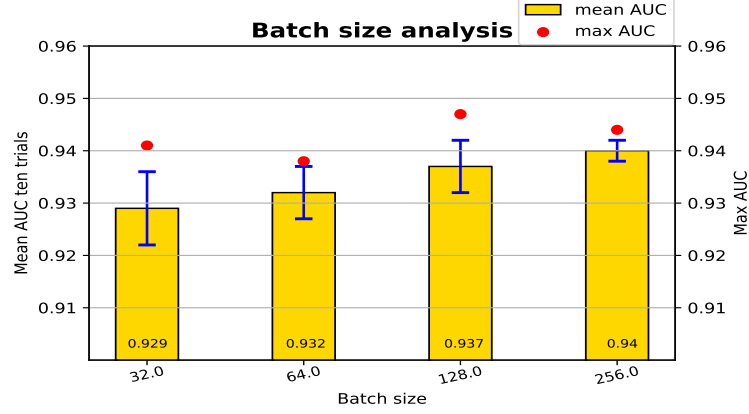


Figure 5.39: Batch size analysis

performed the best, batch size 512 is also evaluated. But it does not train well and the validation accuracy remains stuck near 50%. hence batch size 256 is chosen for the final evaluation.

1613

1614 5.3.12 Total parameter analysis

For determining how well a network work, it is important to find out the total parameter size. For the 'basic network structure' with FC output size 2048 and single FC layer, the total parameters are 3281840. But for network with FC output size 128 has only 1068080 total parameters. FC output size is the main hyperparameter which control the total parameter size for this network.

1620

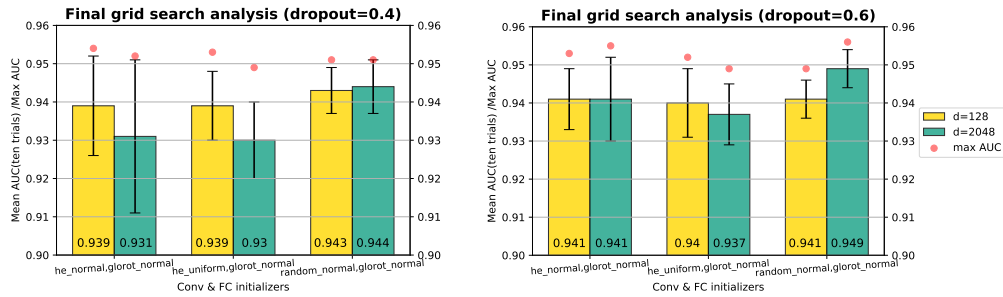
1621 5.3.13 Final grid search

All the best performing hyper parameter values are configured together for this final run. The overall prediction accuracy for the network should be improved with these optimized hyperparameter values. The search space for the final grid search is displayed below in table 5.20.

Ten evaluations are done for each of the test cases, training the network from scratch every time, which is a good test-scenario for the initializers evaluation. The results are displayed in the figure 5.40. The sub-figure 5.40a presents the comparative view of the result obtained by the network with all combination of

Table 5.20: The search space for final grid search for the VGG-Siamese network with Contrastive loss.

Conv filters	Kernel size	FC Layers	FC output	Batch normalization
16	3	1	128, 2048	FALSE
Dropout	Batch size	Optimizer	Initializers	Learning rate
0.4, 0.6	256	Nadam	random normal(Conv) + glorot normal(FC) he uniform(Conv) + glorot normal(FC) he normal(Conv) + glorot normal(FC)	0.0002



(a) Comparative analysis of the final grid search for dropout 0.4.

(b) Comparative analysis of the final grid search for dropout 0.6.

Figure 5.40: Final grid search results for VGG-Siamese network with Contrastive loss. The result for whole search space is visualized here, which is divided in two figures for dropout values 0.4 and 0.6.

configurations mentioned in 5.20, for dropout value 0.4. Similarly for dropout value 0.6 is displayed in the 5.40b.

Discussion

In figure 5.40 in each sub-figure x-axis contains the initializer combinations for Conv and FC layer. For each combination of initializers, two bars are plotted representing mean AUC, obtained for each of the 128 and 2048 FC units value. Y-axis represents mean AUC in ten trials and maximum AUC. The y-axis scales for both sub-figures are matched so that they can be visually compared. The best result for dropout 0.4 is mean AUC (Ten trials) of 0.944 with std of 0.007 and highest AUC value in a single run as 0.95. The best result for dropout 0.6 is mean AUC (Ten trials) of 0.949 with std of 0.005 and highest AUC value in a single run as 0.956. Both aforementioned networks has random normal and glorot normal as Conv and FC layer initializers respectively, this combination is clearly best suited

for the dataset. Also in both cases the FC output size is 2048. Overall highest mean AUC for the VGG-Siamese network with Contrastive loss is 0.949, which is clearly improvement over the state of the art result of 0.894 AUC for score prediction on the same test dataset.

The ROC for the network which yielded maximum AUC for the best configuration, with score of AUC 0.956, is presented in figure 5.41

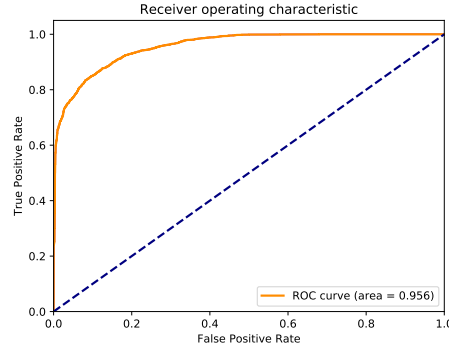


Figure 5.41: ROC for the highest auc recorded for the network with best configuration.

In addition to the hyperparameters, which has best values already determined, presented in table 5.20, rest of the parameter values for the top 4 configurations are displayed in 5.21. The configurations are ordered according to decreasing order for the mean AUC.

Table 5.21: Best hyperparameter values for the Vgg-Siamese with contrastive loss, obtained from the final grid search. In decreasing order of their mean prediction accuracy on the test data. Top 4 results are mentioned.

Config alias	FC output	Dropout probability	Conv initializer	FC initializer	Mean AUC	Std	Max AUC
'Config 1'	2048	0.6	random normal	glorot normal	0.949	0.005	0.956
'Config 2'	2048	0.4	random normal	glorot normal	0.944	0.007	0.951
'Config 3'	128	0.4	random normal	glorot normal	0.943	0.006	0.951
'Config 4'	128	0.6	random normal	glorot normal	0.941	0.005	0.949

With a statistical outlook, five number summary analysis is presented for the top 4 configurations in figure 5.42. it is observed that the median of the configurations are decreasing, similar to the mean AUC. The std for the 'Config 2' is the highest, which is also visible in the box plot, that the data spread is widest

for this configuration. But overall the standard deviation for all 4 results are small and comparable. From the boxplot it is also observable that for 'Config 1' the highest auc value is much higher than other. Not only that, the median lies in the top half of the box, which means at least 5 out of 10 AUC prediction values are higher than 0.95. That makes this configuration very consistent as well.

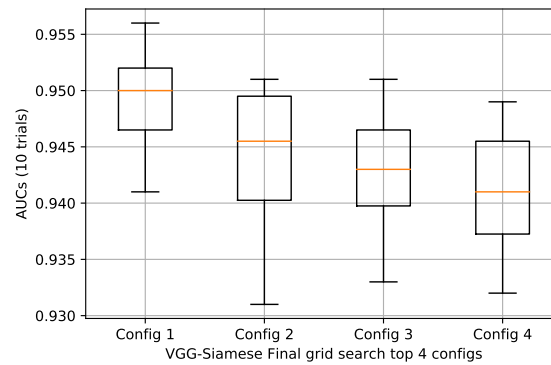


Figure 5.42: Further statistical analysis of top 4 configurations are visualized in box and whisker plots.

With this the hyperparameter optimization is concluded. In the next section the comparative analysis is done for the optimized networks.

5.4 Comparative analysis

So the three network structures that are used, will be compared in this section, in terms of AUC value obtained. Time of execution, total parameters also by using the uncertainty calculated from the Monte Carlo drop out calculations. All our final models are trained with dropouts, as that is the best parameters set up for the network.

5.4.1 AUC comparison

DenseNet two-channel has highest mean AUC(10 Trials) of **0.955**, std 0.009 with max AUC of 0.966. With total parameters of only 51,430. DenseNet-Siamese has highest mean AUC(10 Trials) of **0.921**, std 0.016, Max AUC, 0.95 With total parameters of only 16,725,485. Contrastive loss with VGG-Siamese network have results of mean AUC (Ten trials) of **0.949** with std of 0.005 and highest AUC value in a single run as 0.956. With total parameter size of 3,281,840.

Table 5.22: Comparative analysis on the AUC and total number of parameters in the best performing networks.

Network	Mean AUC	Std	Max AUC	Total params
DenseNet two-channel	0.955	0.009	0.966	51430
DenseNet-Siamese	0.921	0.016	0.95	16725485
VGG-Siamese	0.949	0.005	0.956	3281840

Discussion As seen in table 5.22 the total parameters for the DTC is much lower than the Siamese networks. For both Siamese networks, total parameter size is so large because of the connection of the flattened feature map from each DenseNet branch with fully-connected layer of Siamese branch and following concatenation of the feature maps from both DenseNet branches. If in DenseNet-Siamese each branch has output size P parameters, then after merging or concatenation the total parameters becomes $2 \times P$. If it is connected to the FC layer of output x (For example 2000) then the total parameters involved in that single computation step is $2 \times P \times x$. So in a single step, in Siamese networks the total parameters will increase

1687 by $2 \times$ times P. That is why the total number of parameters for both Siamese are
1688 so much higher than DTC because it does not have this step.

1689

1690 5.4.2 Monte Carlo Dropout analysis

1691 MC Dropout analysis for DenseNet two-channel

1692 In normal training dropout is only applied in the training phase, where it provides
1693 regularization to avoid overfitting. In test time all the connections/nodes remain
1694 present and dropout is not applied, though the weights are adjusted according
1695 to the dropout ratio during training. So every time a prediction on test data
1696 is obtained it is the same, in other words the prediction is deterministic.
1697 There is no randomness. For Monte Carlo dropout the dropout is also applied
1698 in the inference/test time, which introduces randomness, as the connections are
1699 dropped randomly according to the dropout probability. This prediction process
1700 is **stochastic** i.e the model could predict different predictions for same test data.
1701 The main goal of Monte Carlo dropout [15] is to generate random predictions which
1702 can be interpreted as drawing samples from a probabilistic distribution. Monte
1703 Carlo dropout can be compared with approximate inference in Bayesian neural
1704 network.

1705 To enable the dropout at inference time, a previously trained model is chosen
1706 with 0.966 AUC. Out of 10 trials of the best hyperparameter configurations during
1707 the final grid search for DTC 5.2 this model has scored the highest AUC. The
1708 model is trained with 0.2 dropout. For this evaluation the dropout during inference
1709 time is enabled explicitly. 20 stochastic predictions for each of the images are
1710 recorded and the mean prediction and standard deviation is computed which are
1711 mentioned in the 5.43 between two input patches. The more the standard deviation
1712 is, the more uncertain the network is for that specific prediction. Here 20 sequential
1713 images are displayed. The goal is to display how the network prediction works
1714 for matching and non-matching pairs. Here, input pairs from 1000 to 1009 are
1715 matching instances, while 1010 to 1014 are object-object non-matching and rest 5
1716 are object-background non-matching.

1717 Discussion

1718 As seen in 5.43, the maximum std for the matching images are 0.104, for patch
1719 with index 1005, which also has the mean prediction closest to the 0.5, the threshold

in this case. The mean prediction and std for the object-background non-matching are generally low. But it is observed that, even though there is no object in the patch for background, the prediction is not zero. It might be that the network has learned bit of the general noise in sonar patches, which is the only common part here (e.g. patch 1017 with mean prediction 313). But it is very hard to conclude this. However, for object-object non-matching pairs, the network has all the std values higher than 0.1, which means the network has the highest uncertainty for object-object non-matching pairs. This trend was observed through out the whole dataset. Some other images and predictions from same network model with MC dropouts are displayed in figure 5.44 and 5.45. It is observed that the prediction std can be even zero or close to zero for of 20 trials and the mean prediction can be 1 too.

MC Dropout analysis for Siamese networks

Monte Carlo dropout is also applied in for the DenseNet-Siamese and VGG-Siamese network. However, because of the current Keras implementation of dropout, it is probably not applied the same way to both instance of the Siamese branches. Which is causing incorrect weights calculations between the branches and wrong predictions. A simple experiment is conducted with dropout 0.2, 0.4 and 0.6. From the results in table 5.23, with the increase in the dropout ratio, the mean predictions are getting even more erroneous in comparison to the predictions from models without dropout in inference time. No further investigations are conducted for MC dropout application in Siamese network.

Table 5.23: MC dropout prediction example for VGG-Siamese with Contrastive loss for a test image (index 1002). Compared for three models trained with different dropout rates.

Dropout	Normal dropout	Mean (20)	Std (20)	Label(0=matching)
0.2	0.073	1.371	0.175	0
0.4	0.068	1.603	0.183	0
0.6	0.044	2.026	0.126	0

5.4.3 Real prediction analysis

5.4.4 Run time analysis

5.4.5 Ensemble

An ensemble method can be defined as a meta-algorithm which combines several machine learning algorithm into one predictive model in order to improve overall prediction. There are different ensemble techniques such as **Bootstrap aggregating (Bagging) Boosing, stacking, Bayes optimal classifier** etc. These are used with different objectives. The motivation for using ensemble in this thesis came from the following observations. The performance of the DenseNet-Siamese(DS) is good for non-matching pair predictions. DenseNet two-channel(DTC) is overall very good, but most uncertain in object-object non matching pairs. That observation led to the hypothesis that encoding an ensemble of these two classifiers might improve over all predictive capability. For ensemble, the output predictions of both the networks are **averaged** for each test datum. For this test few of the previously trained models of DTC and DS are loaded, and their predictions on the test data are averaged i.e same weightage for DS and DTC both. The ROC AUC calculated on the average prediction is found to be higher than the individual scores each time. Some evaluation results are displayed in table 5.24.

Table 5.24: Ensemble predictions on test data.

DS model AUC	DTC model AUC	Ensemble AUC
0.95	0.959	0.97
0.952	0.959	0.97
0.952	0.963	0.973
0.952	0.966	0.971
0.952	0.972	0.978

Discussion

Ensemble predictions are consistently better than each models separately. If

1765 the underlying models encoding the ensemble has low AUC, the ensemble AUC
1766 is found to be much higher. For example the first result presented in table 5.24
1767 where the ensemble prediction is much higher (0.97 AUC) but the underlying
1768 model predictions (0.95 and 0.959 AUCs). By encoding an ensemble of a DenseNet-
1769 Siamese model with AUC 0.952 and DenseNet two-channel with 0.972 AUC, the
1770 ensemble AUC was **0.978**, which is the highest AUC on test data obtained in any
1771 other tests during this thesis.

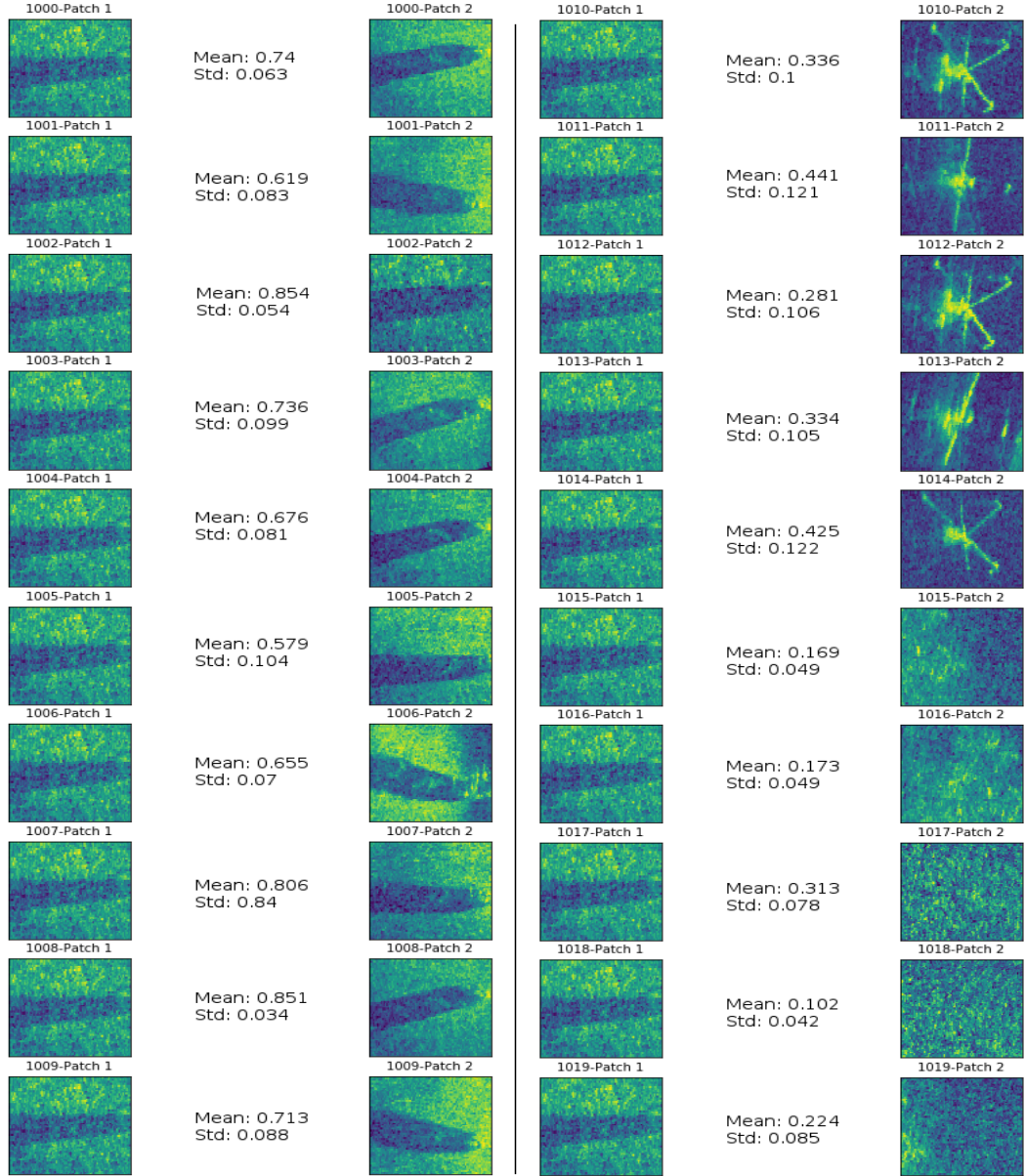


Figure 5.43: 20 pairs of sonar patches, from the test dataset with index 1000 to 1019, and corresponding mean prediction and standard of deviation out of 20 trials with Monte Carlo dropouts.

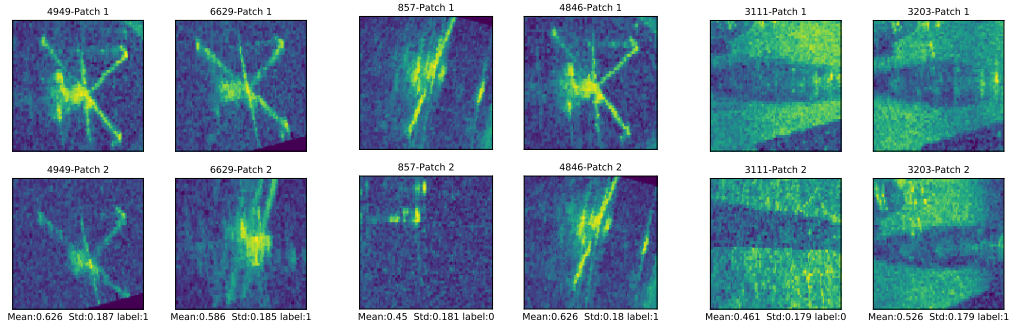


Figure 5.44: Predictions with highest standard deviations, while dropout in inference time is enabled. Mean predictions and corresponding std values (20 trials) and the ground truth (label 1 means matching) are displayed at the bottom of the pairs. It is clear that the low signal-to-noise for sonar is affecting the predictions, and unwanted reflections and occlusions are also adding to the challenge.

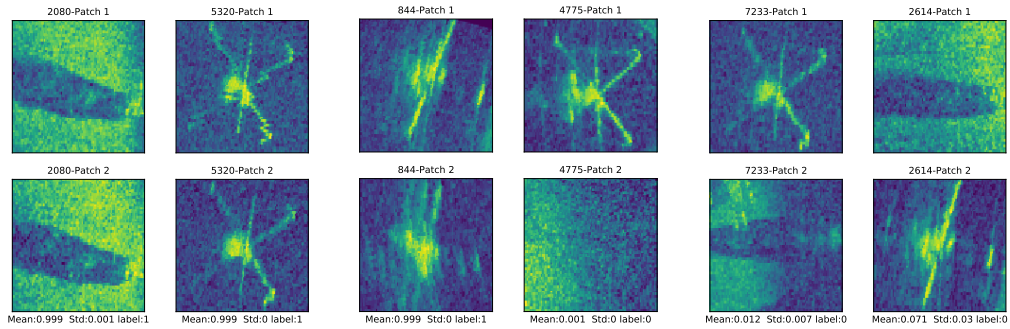


Figure 5.45: Predictions with very low standard deviations, while dropout in inference time is enabled. Highest mean prediction was 1, with std zero or close to zero. Which shows that the network learned some of the similarity functions with great confidence. This is understandable for the test pairs which has very small translational or rotational distance. But index 844 is a very difficult example and prediction with such high confidence is probably unexpected. Again, for object-object non-matching pairs usual std values are much higher than other categories.

Conclusions

6.1 Contributions

1. Hyperparameter optimization and comparative analysis was done for DenseNet two-channel, DenseNet-Siamese and VGG-Siamese network, for sonar data during this work.
2. The state of the art result were improved. TODO explain the best results and state of the art.
3. The use of DenseNet and VGG network as branches of Siamese network, for feature extraction was not applied to sonar data before.
4. Github contribution to Hyperas, which is a Keras wrapper for Hyperopt. It can be used for hyperparameter optimization. But it did not offer any way to evaluate Siamese networks. Initially, Hyperas was chosen as the preferred tool to do hyperparameter optimization. During this stage, a simple program was developed which enables use of Hyperas for Siamese networks. At the Hyperas author's request, this program has been submitted as an example in the Hyperas official repository. Can be found here, https://github.com/maxpumperla/hyperas/blob/master/examples/hyperas_in_intermediate_fns.py.

1791 6.2 Lessons learned

- 1792 • How to encode CNNs to do sonar image processing, without manually de-
1793 signing any features.
- 1794 • It is observed that acoustic vision is very challenging and quite different from
1795 normal optical vision. Which offers scope of a lot of further research. It is
1796 very motivating and challenging at the same time.
- 1797 • Manually doing hyperparameter optimization is a tedious process and needs
1798 some prior experience to be able to predefine the search spaces effectively.

1799 6.3 Future work

- 1800 • FaceNet [35] and triplet loss function is one of the network that was very
1801 interesting for the sonar patch matching.
- 1802 • Do the hyperparameter optimization using a automatic algorithm, for example
1803 Bayesian Optimization [38] or Sequential Model-based Algorithm Configura-
1804 tion (SMAC) [20] or Tree-structured Parzen Estimator (TPE) [1]. Automatic
1805 optimization process involves lesser human involvement and theoretically
1806 converges sooner.

References

- [1] James S Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In *Advances in neural information processing systems*, pages 2546–2554, 2011.
- [2] Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. Signature verification using a” siamese” time delay neural network. In *Advances in neural information processing systems*, pages 737–744, 1994.
- [3] Matthew Brown and David G Lowe. Automatic panoramic image stitching using invariant features. *International journal of computer vision*, 74(1):59–73, 2007.
- [4] François Chollet et al. callbacks keras, December 2018. URL <https://keras.io/callbacks/>.
- [5] François Chollet et al. Conv2D keras, December 2018. URL <https://keras.io/layers/convolutional/>.
- [6] François Chollet et al. Binary cross entropy loss keras, December 2018. URL https://keras.io/losses/#binary_crossentropy.
- [7] François Chollet et al. Initializers keras, November 2018. URL <https://keras.io/initializers/>.
- [8] François Chollet et al. optimizers keras, November 2018. URL <https://keras.io/optimizers/>.
- [9] François Chollet et al. Sigmoid keras, December 2018. URL <https://keras.io/activations/#sigmoid>.

-
- [10] Thibault de Boissiere et al. Densenet implementation, December 2018. URL <https://github.com/tdeboissiere/DeepLearningImplementations/blob/master/DenseNet/densenet.py>.
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. Ieee, 2009.
- [12] Simon Emberton, Lars Chittka, and Andrea Cavallaro. Underwater image and video dehazing with pure haze region segmentation. *Computer Vision and Image Understanding*, 168:145–156, 2018.
- [13] John Hunter et al. (matplotlib team). Make a box and whisker plot, December 2018. URL https://matplotlib.org/api/pyplot_api.html#matplotlib.pyplot.boxplot.
- [14] CNN for visual recognition. Cs231n Stanford University, November 2018. URL <http://cs231n.github.io/neural-networks-3/#update>.
- [15] Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. *arXiv:1506.02142*, 2015.
- [16] The HDF Group. What is hdf5, December 2018. URL <https://support.hdfgroup.org/HDF5/whatishdf5.html>.
- [17] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, volume 1, page 3, 2017.
- [18] Natalia Hurtós, Yvan Petillot, Joaquim Salvi, et al. Fourier-based registrations for two-dimensional forward-looking sonar image mosaicing. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 5298–5305. Ieee, 2012.
- [19] Natalia Hurtós, Narcís Palomeras, Sharad Nagappa, and Joaquim Salvi. Automatic detection of underwater chain links using a forward-looking sonar. In *OCEANS-Bergen, 2013 MTS/IEEE*, pages 1–7. IEEE, 2013.

References

- [20] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *International Conference on Learning and Intelligent Optimization*, pages 507–523. Springer, 2011.
- [21] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [22] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- [23] Juhwan Kim and Son-Cheol Yu. Convolutional neural network-based real-time rov detection using forward-looking sonar image. In *Autonomous Underwater Vehicles (AUV), 2016 IEEE/OES*, pages 396–400. IEEE, 2016.
- [24] K Kim, N Neretti, and N Intrator. Mosaicing of acoustic camera images. *IEEE Proceedings-Radar, Sonar and Navigation*, 152(4):263–270, 2005.
- [25] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [26] Yann LeCun et al. Lenet-5, convolutional neural networks. URL: <http://yann.lecun.com/exdb/lenet>, page 20, 2015.
- [27] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. *arXiv preprint arXiv:1312.4400*, 2013.
- [28] David G Lowe. Object recognition from local scale-invariant features. In *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*, volume 2, pages 1150–1157. Ieee, 1999.
- [29] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.

- [30] Somshubra Majumdar(Github:titu1994) et al. Keras contrib implementation of Densenet, December 2018. URL https://github.com/keras-team/keras-contrib/blob/master/keras_contrib/applications/densenet.py.
- [31] Nicholas Molton, Andrew J Davison, and Ian D Reid. Locally planar patch features for real-time structure from motion. In *Bmvc*, pages 1–10. Citeseer, 2004.
- [32] S Negahdaripour, MD Aykin, and Shayanth Sinnarajah. Dynamic scene analysis and mosaicing of benthic habitats by fs sonar imaging-issues and complexities. In *Proc. OCEANS*, volume 2011, pages 1–7, 2011.
- [33] Minh Tân Pham and Didier Guériot. Guided block-matching for sonar image registration using unsupervised kohonen neural networks. In *Oceans-San Diego, 2013*, pages 1–5. IEEE, 2013.
- [34] Sebastian Ruder. Optimization for deep learning highlights in 2017, November 2018. URL <http://ruder.io/deep-learning-optimization-2017/index.html#understandinggeneralization>.
- [35] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [36] scikit learn. Train test data split, December 2018. URL https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html.
- [37] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *null*, pages 519–528. IEEE, 2006.
- [38] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959, 2012.

References

- [39] Tensorflow. Tensors, December 2018. URL https://www.tensorflow.org/programmers_guide/tensors.
- [40] Matias Valdenegro-Toro. Objectness scoring and detection proposals in forward-looking sonar images with convolutional neural networks. In *IAPR Workshop on Artificial Neural Networks in Pattern Recognition*, pages 209–219. Springer, 2016.
- [41] Matias Valdenegro-Toro. Improving sonar image patch matching via deep learning. In *Mobile Robots (ECMR), 2017 European Conference on*, pages 1–6. IEEE, 2017.
- [42] Peter Vandrish, Andrew Vardy, Dan Walker, and OA Dobre. Side-scan sonar image registration for auv navigation. In *Underwater Technology (UT), 2011 IEEE Symposium on and 2011 Workshop on Scientific Use of Submarine Cables and Related Technologies (SSC)*, pages 1–7. IEEE, 2011.
- [43] Simon Wiesler and Hermann Ney. A convergence analysis of log-linear training. In *Advances in Neural Information Processing Systems*, pages 657–665, 2011.
- [44] Wikipedia. Convolutional neural networks, November 2018. URL https://en.wikipedia.org/wiki/Convolutional_neural_network.
- [45] Wikipedia. Hyperparameter machine learning, November 2018. URL [https://en.wikipedia.org/wiki/Ensemble_averaging_\(machine_learning\)](https://en.wikipedia.org/wiki/Ensemble_averaging_(machine_learning)).
- [46] Bangpeng Yao, Gary Bradski, and Li Fei-Fei. A codebook-free and annotation-free approach for fine-grained image categorization. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 3466–3473. IEEE, 2012.
- [47] Sergey Zagoruyko and Nikos Komodakis. Learning to compare image patches via convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4353–4361, 2015.
- [48] Jure Zbontar and Yann LeCun. Stereo matching by training a convolutional neural network to compare image patches. *Journal of Machine Learning Research*, 17(1-32):2, 2016.