

## Part I

Q1 The minimum information required is  $I_n$  in Gaussian pyramid.

By definition of Laplacian,  $F$  denotes upsampling

$$L_i = I_i - F(I_{i+1}) \Rightarrow I_k = F(I_{i+1}) + L_i$$

$$\text{Base case } I_0 = F(I_1) + L_0$$

$$= F(F(I_2) + L_1) + L_0$$

{

$$= F(F(\dots(F(I_n) + L_{n-1}) + \dots + L_1) + L_0$$

$$= \sum_{i=1}^n F(I_i + L_{i-1})$$

$\therefore$  the minimum information require is  $I_n$ .

Q2. Let  $x$  be the input,  $W$  be the layer weight, and  $F$  be the activation function matrix of function  $f$ .

$$f(Wx + b) = FWx + Fb$$

So it is still a linear relationship for the activation function layers.

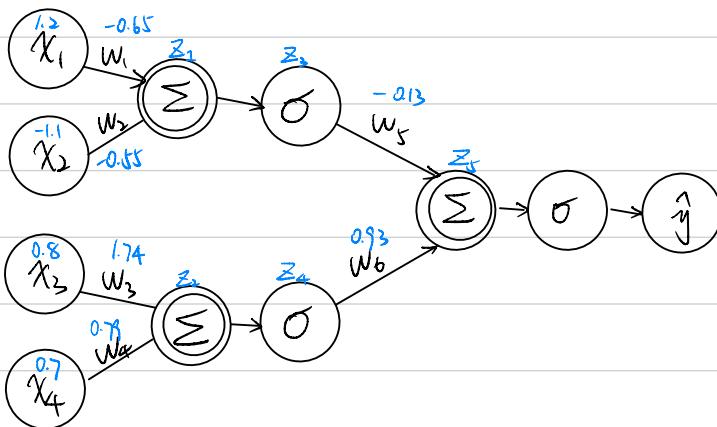
$$\text{For } n \text{ layers we have } \hat{y} = f(W_n f(W_{n-1} \dots f(W_1 x))) = (f \circ W_n) \circ \dots \circ (f \circ W_1) x$$

which is also a linear relationship can be rewritten as  $Wx$ .

$\therefore$  The number of layers has no impact on the network.

Q3.

(a)



(b) Value marked in blue are for part (b) for easy visualization.

$$\Sigma_1 = w_1 x_1 + w_2 x_2 = -0.1750 \quad \Sigma_2 = w_3 x_3 + w_4 x_4 = 1.9450$$

$$\Sigma_3 = O(\Sigma_1) = 0.4564 \quad \Sigma_4 = O(\Sigma_2) = 0.8749$$

$$\Sigma_5 = w_5 \Sigma_3 + w_6 \Sigma_4 = 0.7543 \quad \hat{y} = O(\Sigma_5) = 0.6801$$

$$\frac{\partial L}{\partial \hat{y}} = -2(y - \hat{y}) = -2(1.0 - 0.6801) = -0.6398$$

$$\frac{\partial L}{\partial z_5} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_5} = \frac{\partial L}{\partial \hat{y}} \cdot (1 - \hat{y}) = -0.6398 \cdot (0.6801 \cdot (1 - 0.6801)) = -0.1392$$

$$\frac{\partial L}{\partial z_4} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_4} = \frac{\partial L}{\partial \hat{y}} \cdot w_6 = -0.1392 \cdot 0.93 = -0.1295$$

$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} = \frac{\partial L}{\partial \hat{y}} \cdot (\Sigma_4(1 - \Sigma_4)) = -0.1295 \cdot (0.8749 \cdot (1 - 0.8749)) = -0.0142$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_3} = \frac{\partial L}{\partial \hat{y}} \cdot X_3 = -0.0142 \cdot 0.8 = -0.0114$$

Q4: For each filter, the size is  $4 \times 4 \times 50 = 800$

The stride is 2, so there are  $\frac{12+2-1-4}{2} + 1 = 6$  positions and 20 filters

For each entry we have  $4 \times 4 = 16$  multiplication and  $4 \times 4 - 1 = 15$

additions  $20 \times 6 \times 6 \times 50 \times (16 + 15) = 1116000$  FLOPs without bias

and  $20 \times 6 \times 6 \times 50 \times (16 + 15 + 1) = 1152000$  FLOPs with bias

The maxpooling for  $6 \times 6 \times 20$  size require  $\frac{6-3}{1} + 1 = 4$ , meaning max pooling

size is  $4 \times 4$ . and each max pooling takes  $3 \times 3 - 1 = 8$  FLOPs.

$\therefore 4 \times 4 \times 20 \times 8 = 2560$  FLOPs

So the number of FLOPs is  $1116000 + 2560 = 1118560$  FLOPs without bias

and  $1152000 + 2560 = 1154560$  FLOPs with bias

Q5

For C1:  $32 - 28 + 1 = 5$  kernel size

$$(5 \times 5 + 1) \times 6 = 156 \text{ parameters.}$$

S2: 0 parameters

C3:  $14 - 10 + 1 = 5$  kernel size

$$(5 \times 5 \times 6 + 1) \times 16 = 2416 \text{ parameters.}$$

S4: 0 parameters

C5:  $5 \times 5 \times 16 \times 120 = 48000$  parameters.

F6:  $120 \times 84 + 84 \times 10 = 10920$  parameters.

Total:  $156 + 0 + 2416 + 0 + 48000 + 10920 = 61492$  parameters.

Q6: Assume input  $x$ , weight  $w$  and bias  $b$ .

$$\therefore L = \sigma(wx + b)$$

$$\frac{\partial L}{\partial x} = w\sigma(wx + b)(1 - \sigma(wx + b)) \\ = wL(1 - L)$$

From the result, we can see we only need  $L$  and  $w$  in order to compute the derivative of output, no need for the inputs.

Q7. (a) The output range of the function is  $(-1, 1)$  because  $\lim_{x \rightarrow -\infty} = 1$  and  $\lim_{x \rightarrow \infty} = -1$   
 and the output range of logistic function is  $(0, 1)$

$$\begin{aligned}(b) \frac{d}{dx} \tanh(x) &= \frac{2e^{-2x}(1+e^{-2x}) + 2e^{-2x}(-e^{-2x})}{(1+e^{-2x})^2} \\ &= \frac{4e^{-2x}}{(1+e^{-2x})^2} \\ &= 1 - \tanh^2(x)\end{aligned}$$

$$\text{where } \tanh(x) = \frac{1-e^{-2x}}{1+e^{-2x}} = 2\sigma(2x) - 1$$

$\therefore$  gradient of tanh can be formulated as a function of logistic function.

(c) The usage of each activation function depends on the nature of the output range.

For example, the  $\tanh(x)$  preserved the sign of input, it may be more useful in binary classification.

And the logistic activation can be used for probability predication, since the output range is  $(0, 1)$

## Part II

### Task I:

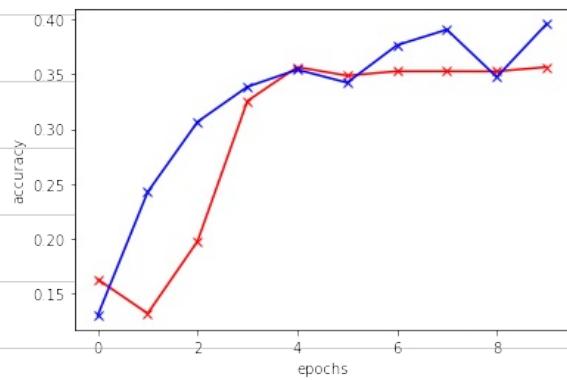
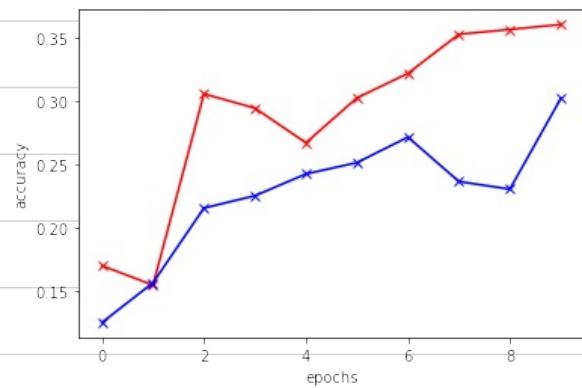
DBI: The background is very simple or even removed, we can clearly identify the main object is the dog in the picture.

SDD: The sources of images seem to vary. The background is not so clear and may even hard to separate the dog. There are a lot of images contains non-relevance object like human, and decoration on the dog, food... And some of the images even contain water mark.

## Task II

with dropout

Without dropout



red : test accuracy

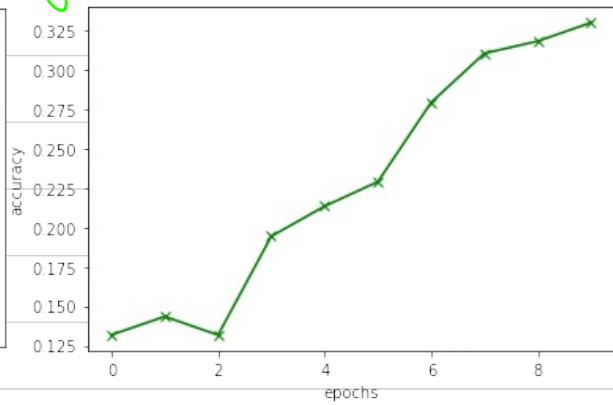
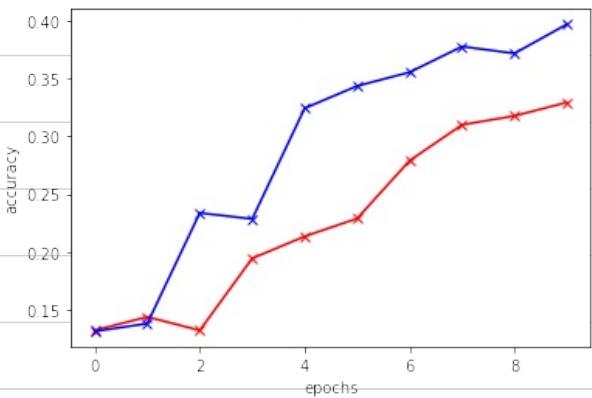
blue : train accuracy

From two reports, we can see there is no obvious increase in the overall accuracy, but the fit speed seems to be faster without dropout. It may because the model we are training is relatively small and dropout in the small data made it even hard to fit.

### Task III

a. ResNet-18 on DBI

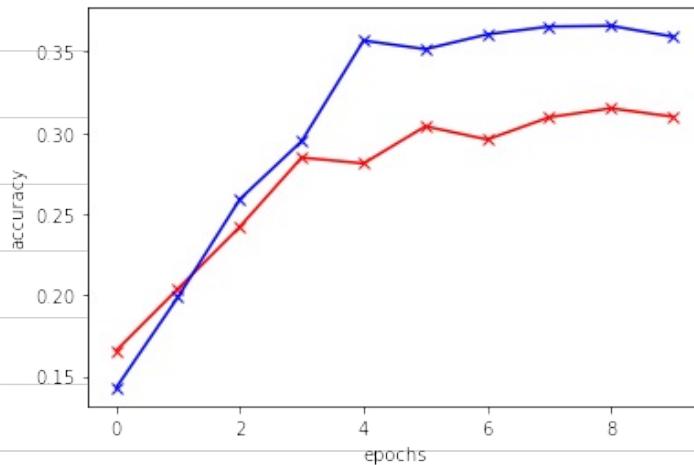
red: test blue: train green: validation



Comparing with CNN model, we can see the overall accuracy is about the same with test accuracy lower than the CNN model. But the trend seems to keep growing. So I think this model can achieve better accuracy than CNN with more training.

## b. ResNet-18 on SDD

red: test blue: train

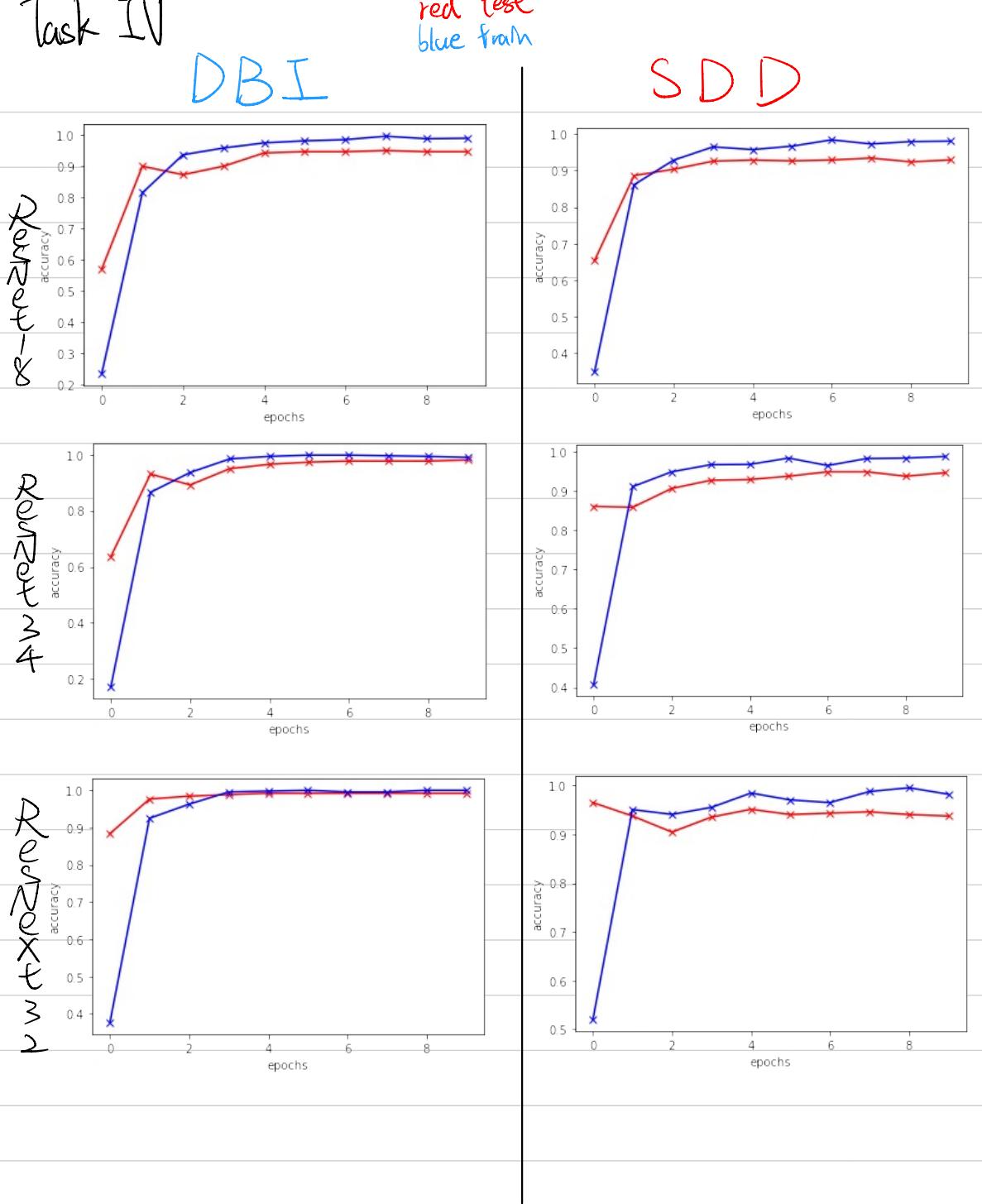


The overall accuracy is a little bit lower than the dataset on DB1. The reason may be the background is more complex in SDD and the model may not correctly learned what is dog. It may learned the background like human as dog.

# Task IV

red test  
blue train

DBI

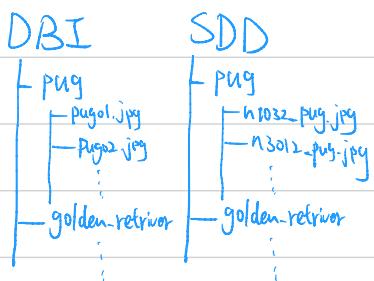


From the result, we can see these networks achieved similar results at the end with over 0.95 accuracy. In general, all networks are performing better in DBI set which is expected because the set is more clean. Among all, the ResNeXt32 seems to win from the rest in the DBI set. There is no clear winner in SDD set but ResNeXt32 still has a little bit better accuracy among all.

## Task V

In order to classifying the image is from SDD or DBI. I reorganized the datasets as following.

Before



After



I moved all images in the sub folder to the parent folder and make a dataset with

subset DBI and SDD. All images under DBI will have label 'DBI' and all images under SDD will have label 'SDD'. I modified the best resulting network from previous task - ResNet32 to get it learn from 2 features. The result is as following.

