

A.



UNIVERSITY OF  
**WEST LONDON**  
The **Career** University

**Deep Learning (CP70071E\_05-FEB-  
24\_31-MAY-24\_A\_SEM2)**

**Topic – Log Book**

**Student ID: [21600713]**

**Full Name: [Mohammed Haris Shaikh]**

## II. Table of Contents

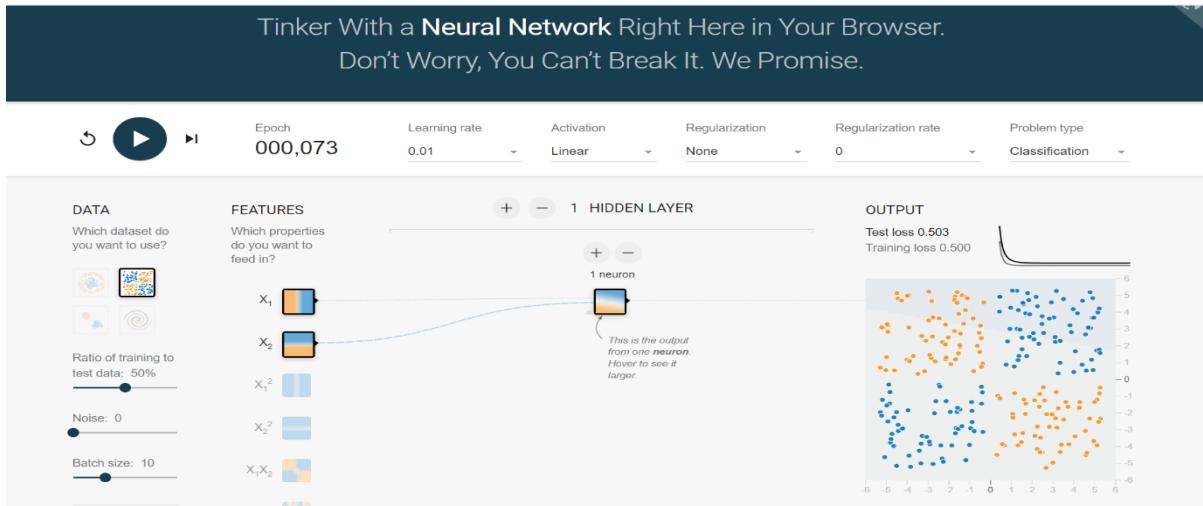
A. ....	0
<b>III. WEEK 1</b> .....	4
A. Task 1: The model as given combines our two input features into a single neuron. Will this model learn any non-linearities? Run it to confirm your guess. ....	4
B. Task 2: Try increasing the number of neurons in the hidden layer from 1 to 2, and also try changing from a Linear activation to a nonlinear activation like ReLU. Can you create a model that can learn nonlinearities? Can it model the data effectively? .....	4
C. Task 3: Try increasing the number of neurons in the hidden layer from 2 to 3, using a nonlinear activation like ReLU. Can it model the data effectively? How model quality vary from run to run? .....	5
D. Task 4: Continue experimenting by adding or removing hidden layers and neurons per layer. Also feel free to change learning rates, regularization, and other learning settings. What is the smallest number of neurons and layers you can use that gives test oss of 0.177 or lower?.....	5
E. Task 5: Does increasing the model size improve the fit, or how quickly it converges? Does this change how often it converges to a good model? For example, try the following architecture: • First hidden layer with 3 neurons. • Second hidden layer with 3 nurons. • Third hidden layer with 2 neurons.....	6
<b>IV. WEEK 2</b> .....	7
A. Resource: With assumption that you have access to Microsoft Excel, create a new template with the following setup. From there, you start implement the perceptron binary classification for the AND function approximation based on the formulas discussed in the lecture. .....	7
B. Expected outcome: You are required to plot the Cost function value at every epoch, up to 70 epochs. Your plot should be inline with the Figure 3. ....	7
C. Investigate if the Cost function can reach less than 0.09. If it is possible, by how many epoch the Cost function reaches less than 0.09? .....	8
<b>V. WEEK 3</b> .....	9
A. 1. Please use the followimg link for Hands-on session: .....	9
B. Workshop 2a .....	13
C. Workshop 2b .....	15
1. Week 3 .....	17
<b>VI. WEEK 4</b> .....	20
A. Resource: With assumption that you have access to Microsoft Excel, create a new template with the following setup. From there, you start implement the single-	

layer perceptron multiclass classification for the Binary to Decimal decoder function approximation based on the formulas discussed in the lecture. ....	20
1. Investigate if we should put an additional Sigmoid function before the Softmax function. Explain your finding along with your analysis evidence. ....	21
2. Impact of Activation Function Choices in Multi-Class Classification .....	21
B. Please run the following Colab notebook and report your results in your logbook: 23	
<b>VII. WEEK 5 .....</b>	<b>27</b>
A. Resource: With assumption that you have access to Microsoft Excel, create a new template with the following setup. From there, you start implement the multilayer perceptron logistic regression for solving XOR function approximation based on the formulas discussed in the lecture. ....	27
1. Expected outcome: You are required to plot the Cost function value at every epoch, up to 70 epochs. Your plot should be inline with the Figure 3.....	27
2. Replace Sigmoid function with Tanh function. Explain your finding along with your analysis evidence.....	28
B. Please complete the following notebook and the attached workshop, record the outcome in your log book and demonstrate your understand to your best level. ....	29
<b>VIII. WEEK 6 .....</b>	<b>31</b>
A. Please complete the following notebook, record the outcome in your log book and demonstrate your understand to your best level. ....	31
<b>IX. WEEK 7 .....</b>	<b>39</b>
A. Please complete the following notebook, record the outcome in your log book and demonstrate your understand to your best level. ....	39
<b>X. Week 8.....</b>	<b>46</b>
A. Please complete the following notebook, record the outcome in your log book and demonstrate your understand to your best level. ....	46
<b>XI. Week 9.....</b>	<b>53</b>
A.  Week9 - CNN Models(1).ipynb .....	53
B. Please complete the following notebook and the attached workshop, record the outcome in your log book and demonstrate your understand to your best level. ....	55
1. Analyze Computational Performance of AlexNet:.....	63
2. Simplified AlexNet for Fashion-MNIST Classification.....	65
3. Modified Alexnet for 28x28 images.....	67
<b>XII. Week 10 and 11.....</b>	<b>68</b>

- A. Please complete the workshop and the following notebook, record the outcome in your log book and demonstrate your understand to your best level. ..... 68
- B. Analyzing Textual Time Travel Data with RNN Models..... 72**

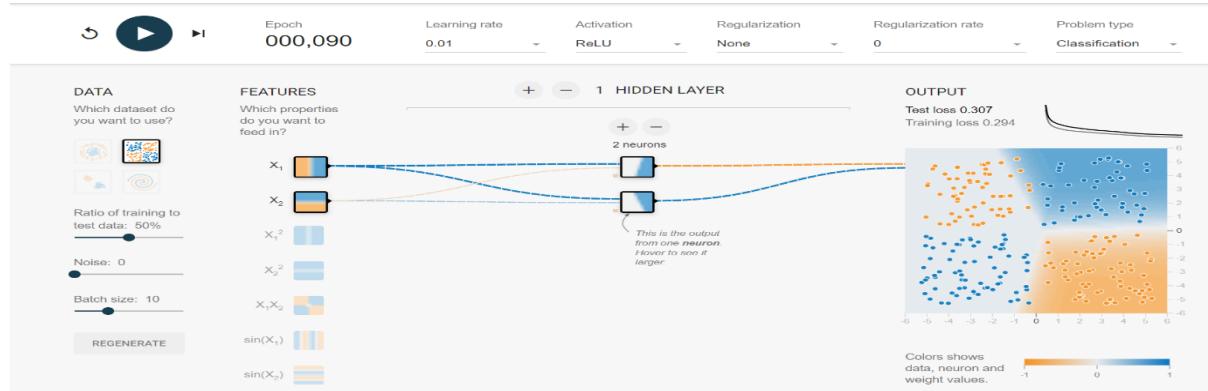
### III. WEEK 1

A. Task 1: The model as given combines our two input features into a single neuron. Will this model learn any non-linearities? Run it to confirm your guess.



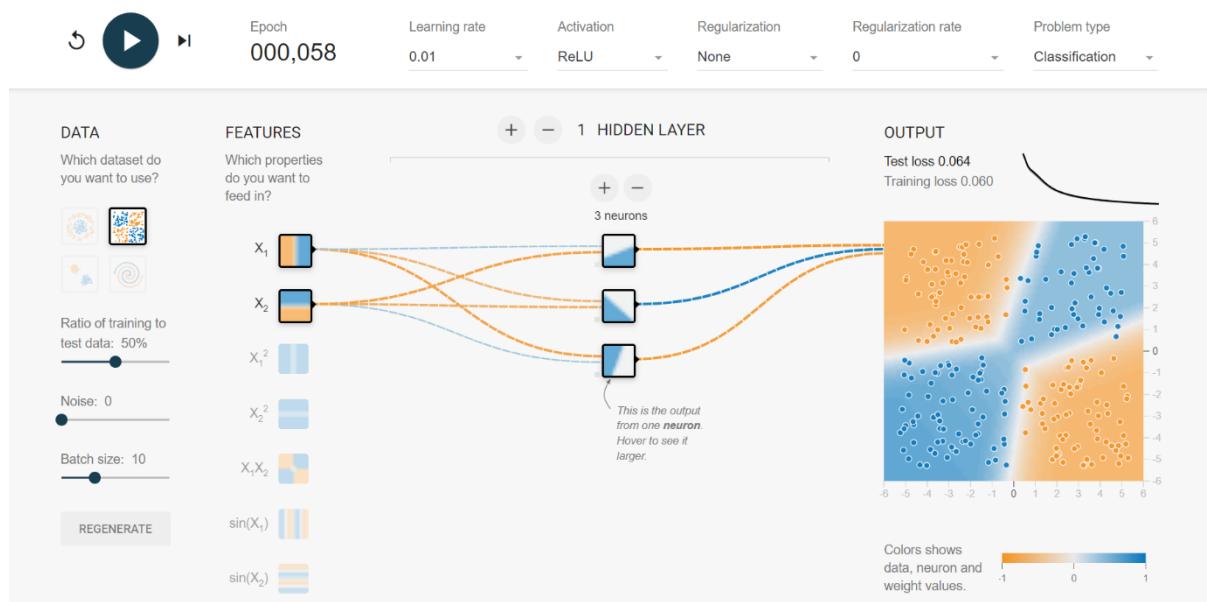
**Ans-**This model with a single neuron and a linear activation function will not be able to learn any nonlinearities between the two input features.

B. Task 2: Try increasing the number of neurons in the hidden layer from 1 to 2, and also try changing from a Linear activation to a nonlinear activation like ReLU. Can you create a model that can learn nonlinearities? Can it model the



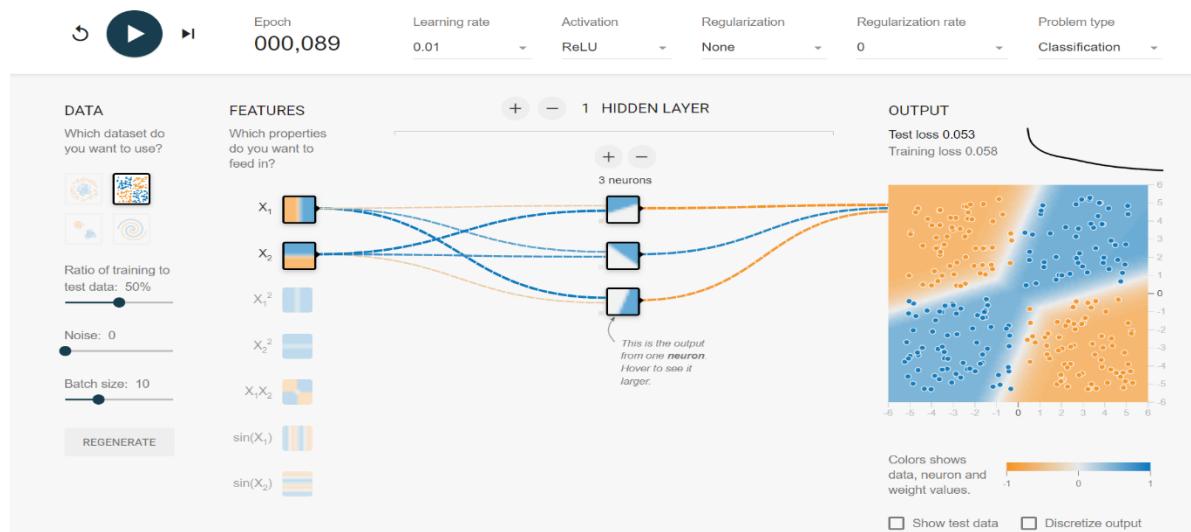
**Ans-** Increasing the number of neurons in the hidden layer and changing the activation function from Linear to ReLU can indeed help a neural network learn nonlinearities and model the data more effectively. However above model with only one neuron added will not have any effect still model cannot predict data efficiently .

C. Task 3: Try increasing the number of neurons in the hidden layer from 2 to 3, using a nonlinear activation like ReLU. Can it model the data effectively? How model quality vary from run to run?



**Ans-Above model with three neurons and using relu it predicted data effectively, unfortunately if we run again and again data effectiveness decreases sometime it perfect sometimes having test loss 0.25.**

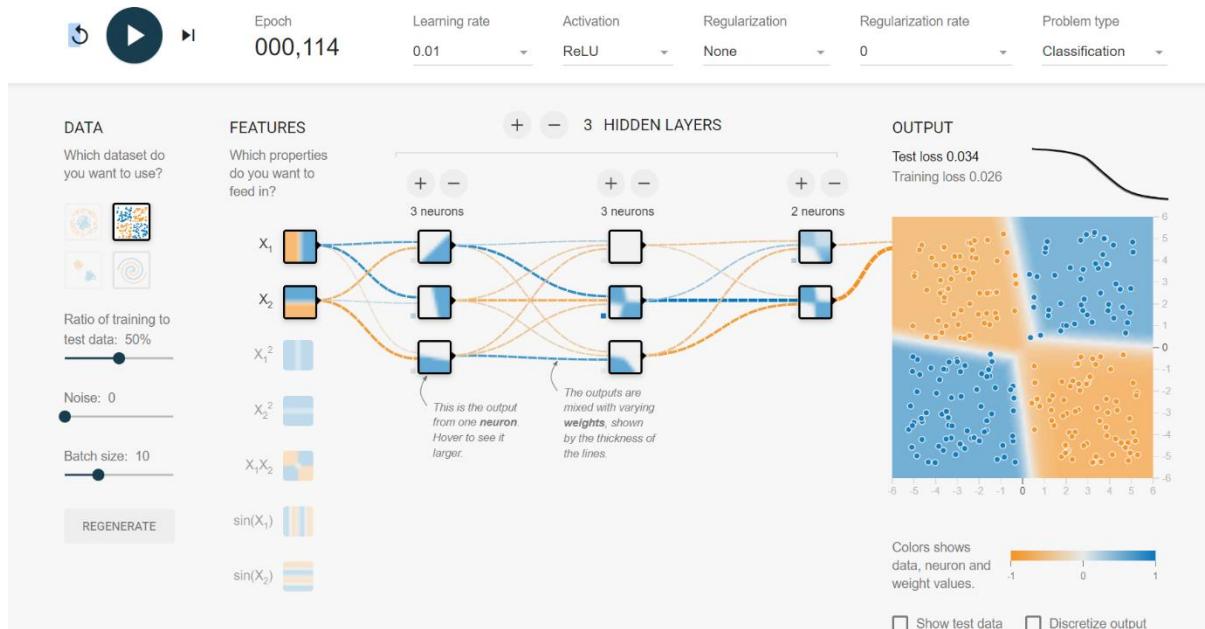
D. Task 4: Continue experimenting by adding or removing hidden layers and neurons per layer. Also feel free to change learning rates, regularization, and other learning settings. What is the smallest number of neurons and layers you can use that gives test loss of 0.177 or lower?



**Ans- I played with parameter and explored models couple times ,However above model with 1 hidden layer and 3 neurons using relu and learning rate - 0.01 I get model with less than 0.177 test loss.**

E. Task 5: Does increasing the model size improve the fit, or how quickly it converges? Does this change how often it converges to a good model? For example, try the following architecture:

- First hidden layer with 3 neurons.
- Second hidden layer with 3 neurons.
- Third hidden layer with 2 neurons



**Ans- Definitely with increasing hidden layers and neuron effects data quality and help to predict non-linearity but convergence to global minima required more epochs in my case, Sometime it get overfit the model, However while running model again and again more probably get good results as compare to low hidden layers and neurons model.**

## IV. WEEK 2

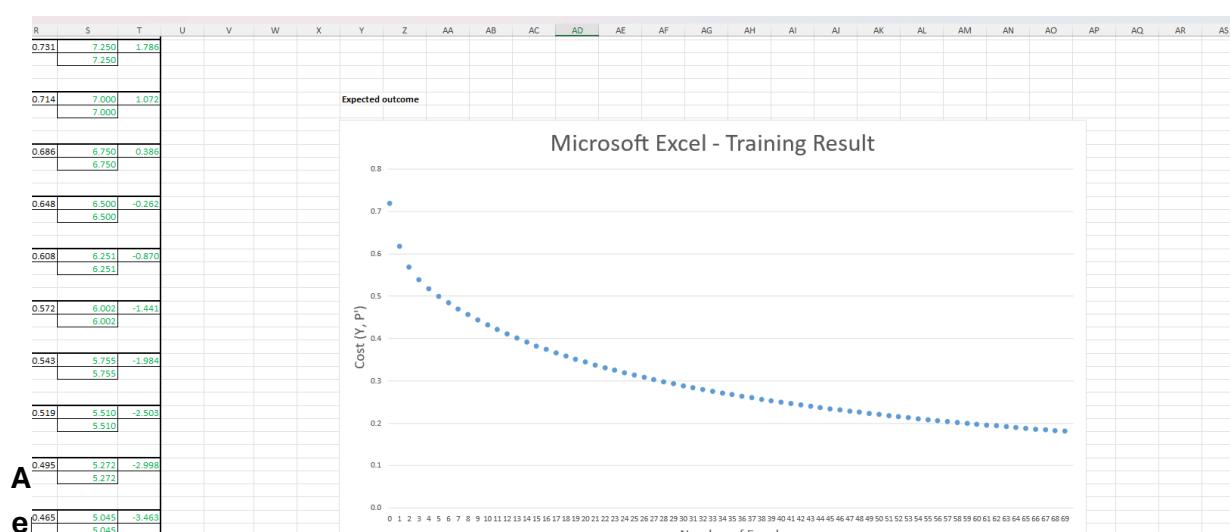
**A. Resource:** With assumption that you have access to Microsoft Excel, create a new template with the following setup. From there, you start implement the perceptron binary classification for the AND function approximation based on the formulas discussed in the lecture.

Epoch	X	Y	W	b	L	Rate	Z	p	$y \log(p)$	$(1-y) \log(1-p)$	Cost(Y, P)	$dC/dp$	$dP/dZ$	$dC/dZ$	$dZ/db$	$dC/dW$	$dC/db$	W'updated weight	b'updated bias for next epoch
0	0	0	0	10	10	1	1.000	1.000	-	-	12.50001	5500.000	0.000	0.250	1.000	0.250	0.750	9.250	
	0	1	0	10	10	1	1.000	0.000	20.000	20.000	121291294.749	0.000	0.250	0.000	0.250	0.000	0.750	9.250	
	1	0	0	10	10	1	20.000	0.000	20.000	20.000	121291294.749	0.000	0.250	0.000	0.250	0.000	0.750	9.250	
	1	1	1	10	10	1	30.000	0.000	0.000	0.000	121291294.749	0.000	0.250	0.000	0.000	0.000	0.750	9.250	
1	0	0	0	9.250	9.250	1	9.250011	1.000	0.000	9.250	11.81253	2601.421	0.000	0.250	1.000	0.250	0.750	9.500	8.500
	0	1	0	9.250	9.250	1	19.00001	1.000	0.000	19.000	44621080.906	0.000	0.250	0.000	0.250	0.000	0.750	9.500	
	1	0	0	9.250	9.250	1	19.00001	1.000	0.000	19.000	44621080.906	0.000	0.250	0.000	0.250	0.000	0.750	9.500	
	1	1	1	9.250	9.250	1	28.75001	1.000	0.000	0.000	-0.250	0.000	0.000	0.000	0.000	0.000	0.750	9.500	
2	0	0	0	9.500	8.500	1	8.500035	1.000	0.000	8.500	11.12508	1228.886	0.000	0.250	1.000	0.250	0.750	9.250	7.750
	0	1	0	9.500	8.500	1	18.00004	1.000	0.000	18.000	16415573.402	0.000	0.250	0.000	0.250	0.000	0.750	9.250	
	1	0	0	9.500	8.500	1	18.00004	1.000	0.000	18.000	16415573.402	0.000	0.250	0.000	0.250	0.000	0.750	9.250	
	1	1	1	9.500	8.500	1	27.50004	1.000	0.000	0.000	-0.250	0.000	0.000	0.000	0.000	0.000	0.750	9.250	
3	0	0	0	9.250	7.750	1	7.750086	1.000	0.000	7.751	10.43767	580.693	0.000	0.250	1.000	0.250	0.750	9.000	7.000
	0	1	0	9.250	7.750	1	17.00001	1.000	0.000	17.000	6039259.275	0.000	0.250	0.000	0.250	0.000	0.750	9.000	
	1	0	0	9.250	7.750	1	17.00001	1.000	0.000	17.000	6039259.275	0.000	0.250	0.000	0.250	0.000	0.750	9.000	
	1	1	1	9.250	7.750	1	26.25009	1.000	0.000	0.000	-0.250	0.000	0.000	0.000	0.000	0.000	0.750	9.000	
4	0	0	0	9.000	7.000	1	7.000194	0.999	0.000	7.001	9.75037	274.461	0.001	0.250	1.000	0.250	0.750	8.750	6.250
	0	1	0	9.000	7.000	1	16.00019	1.000	0.000	16.000	2221958.695	0.000	0.250	0.000	0.250	0.000	0.750	8.750	
	1	0	0	9.000	7.000	1	16.00019	1.000	0.000	16.000	2221958.695	0.000	0.250	0.000	0.250	0.000	0.750	8.750	
	1	1	1	9.000	7.000	1	25.00019	1.000	0.000	0.000	-0.250	0.000	0.000	0.000	0.000	0.000	0.750	8.750	
5	0	0	0	8.750	6.250	1	6.250422	0.998	0.000	6.252	9.06330	129.808	0.002	0.250	1.000	0.250	0.750	8.500	5.500
	0	1	0	8.750	6.250	1	15.00042	1.000	0.000	15.000	817599.311	0.000	0.250	0.000	0.250	0.000	0.750	8.500	
	1	0	0	8.750	6.250	1	15.00042	1.000	0.000	15.000	817599.311	0.000	0.250	0.000	0.250	0.000	0.750	8.500	
	1	1	1	8.750	6.250	1	23.75042	1.000	0.000	0.000	-0.250	0.000	0.000	0.000	0.000	0.000	0.750	8.500	
6	0	0	0	8.500	5.501	1	5.500903	0.996	0.000	5.505	8.37670	61.478	0.004	0.249	1.000	0.250	0.749	8.250	4.752
	0	1	0	8.500	5.501	1	14.00095	1.000	0.000	14.001	300923.058	0.000	0.250	0.000	0.250	0.000	0.750	8.250	
	1	0	0	8.500	5.501	1	22.50095	1.000	0.000	0.000	-0.250	0.000	0.000	0.000	0.000	0.000	0.750	8.250	
	1	1	1	8.500	5.501	1	31.25095	0.991	0.000	4.743	7.69109	29.000	0.008	0.248	1.000	0.250	0.748	8.000	4.004
7	0	0	0	8.250	4.752	1	4.75192	0.991	0.000	4.752	110816.234	0.000	0.250	0.000	0.250	0.000	0.750	8.000	
	0	1	0	8.250	4.752	1	13.00192	1.000	0.000	13.002	110816.234	0.000	0.250	0.000	0.250	0.000	0.750	8.000	
	1	0	0	8.250	4.752	1	13.00192	1.000	0.000	13.002	110816.234	0.000	0.250	0.000	0.250	0.000	0.750	8.000	
	1	1	1	8.250	4.752	1	21.25192	1.000	0.000	0.000	-0.250	0.000	0.000	0.000	0.000	0.000	0.750	8.000	
A	0	0	0	8.000	4.004	1	4.004062	0.982	0.000	4.022	7.00757	13.955	0.018	0.246	1.000	0.250	0.746	7.750	3.250
	0	1	0	8.000	4.004	1	12.00406	1.000	0.000	12.004	40854.588	0.000	0.250	0.000	0.250	0.000	0.750	7.750	
	1	0	0	8.000	4.004	1	12.00406	1.000	0.000	12.004	40854.588	0.000	0.250	0.000	0.250	0.000	0.750	7.750	
	1	1	1	8.000	4.004	1	20.00406	1.000	0.000	0.000	-0.250	0.000	0.000	0.000	0.000	0.000	0.750	7.750	

**Ans- As per given information all details filled and applied formula for finding cost function and back propagation(updated weight and bias for next epoch)**

You will find excel file here

**B. Expected outcome:** You are required to plot the Cost function value at every epoch, up to 70 epochs. Your plot should be inline with the Figure 3.



C. Investigate if the Cost function can reach less than 0.09. If it is possible, by how many epoch the Cost function reaches less than 0.09?

**Ans- As we can see after 37 epochs cost function reaches 0.099 and still decreasing after every epochs, in my opinion if it still continue this model gonna be overfit which is not at all good .**

## V. WEEK 3

### A. 1. Please use the following link for Hands-on session:

<https://colab.research.google.com/drive/1wuJFhZ0DGvAu1YjWCImFPkc6I3FWdLwq?usp=sharing>

### Binary Sentiment Classification of IMDB Reviews

This project investigates sentiment analysis for movie reviews using a deep learning model. The objective is to classify reviews as positive or negative.

#### Data and Preprocessing:

- **Dataset:** The IMDB movie review dataset was used for this binary classification task.
- **Preparation:** A subset of 25,000 reviews was extracted and preprocessed for the model. This likely involved steps like text cleaning, tokenization, and vectorization to convert text data into a format suitable for machine learning algorithms.

The IMDB dataset contains 50000 reviews from the Internet Movie Database (imdb.com). By default, the data is split into 25000 reviews for training and 25000 reviews for testing and each split set is balanced, containing 50% negative reviews and 50% positive reviews.

This dataset is packaged with Keras and has already been pre-processed for you, i.e., the words have already been converted into sequences of integers as we will show below.

We specify num\_words = 10000, this means that we only want to keep the top 10000 most frequently occurring words in the training data and rare words would be discarded. The reason for this is so we have a fixed vector of 10000 going into our neural network later on, and also to keep things at a manageable size.

```
[1]: from tensorflow.keras.datasets import imdb
      (train_data, train_labels), (test_data, test_labels) = imdb.load_data(num_words=10000)

[2]: Downloading data from https://storage.googleapis.com/tensorflow/tf_keras/datasets/imdb.npz
      1/642/89/1/464/89 [=====] - 0s/step
```

Lets inspect the data.

Executing the cell below, we see that the total number of records in the training dataset is 25000 rows.

We can also see that the first movie review at index 0 is a sentence containing 218 words.

The actual data in the first record (at index 0) contains a sentence where each word has been converted to an integer. This integer corresponds to a word in a word dictionary.

```
[3]: print(len(train_data))
      print(len(train_data[0]))
      print(train_data[0])

[4]: 25000
      218
      [1, 14, 22, 16, 43, 539, 973, 1622, 1385, 65, 458, 468, 66, 2941, 4, 173, 36, 216, 5, 25, 59, 45, 838, 112, 50, 679, 2, 9, 35, 489, 284, 5, 156, 4,
```

0 train\_labels[0]

1

The cell below shows that the length of the review with the most words is 999.

```
[5]: max([len(sequence) for sequence in train_data])
```

999

Deoding reviews back to text

Convert the word indices of the review at index 0 of the training dataset back to a string sentence. Now we can read the review. We will replace missing words with a question mark. Remember we only kept the most frequently used words and discarded the rest.

```
[6]: word_index = imdb.get_word_index()
      reverse_word_index = dict([(value, key) for (key, value) in word_index.items()])
      decoded_sentence = ' '.join([reverse_word_index.get(i - 3, '?') for i in train_data[0]])
      decoded_sentence
```

Downloading data from https://storage.googleapis.com/tensorflow/tf\_keras/datasets/imdb\_word\_index.json  
16422/16422 [=====] - 0s/step

? this film was just brilliant casting location scenery story direction everyone's really suited the part they played and you could just imagine being there except ? is an amazing actor now the same being director ? father came from the same scottish island as myself so i loved the fact there was a real connection with this film the slyty rebeccah throughout the film were at it was just brilliant so much that i bought the film as soon as it was released for ? and would recommend it to everyone to watch and the fly fishing was amazing really cruddy at the end so sad you know what they say if you cry at a film it must have been good and this definitely was also ? he two little boy's that played the ? of momma and pappa they were just brilliant children are often left out of the ? list i think because the stars that play them all grown up are such a big profile for the whole film but these children are amazing and should be praised for what they have done don't you thi ?

Preparing the data

Encoding the integer sequences via multi-hot encoding

You cannot pass the training data just yet to the neural network. Each review has different lengths. Lets turn them into a vector of zeros and ones and convert the sequence into a 10000-dimensional vector.

For example, if a sequence/review contained the words 'Very good' and its integer sequence in our training dataset was [10, 15] (10=very and 15=good). When we convert this sequence to a 10000 dimensional vector, all would be 0's except for indices 10 and 15.

```
[46]: import numpy as np
      def vectorize_sequences(sequences, dimension=10000):
          results = np.zeros((len(sequences), dimension))
          for i, sequence in enumerate(sequences):
              for j in sequence:
                  results[i, j] = 1.
          return results
      x_train = vectorize_sequences(train_data)
      x_test = vectorize_sequences(test_data)
```

We store these in new variables x\_train and x\_test. In the cell below, we see that the review at position 0 has now been multi-hot encoded. Compare this with the same review in the old variable train\_data[0] from above.

```
[48]: print(x_train[0])
      [0. 1. 1. ... 0. 0. 0.]
```

x\_train refers to our training data, y\_train contains the labels for it. Similarly, x\_test refers to our testing data and y\_test contains the labels for it.

## Model Building and Overfitting:

- Architecture (Initial):** A deep learning model with three dense layers was constructed:
  - First Layer: 16 neurons
  - Second Layer: 16 neurons
  - Third Layer: 1 neuron with sigmoid activation (common for binary classification)
- Training:** The model was trained for 20 epochs using separate training and validation sets (`x_train`, `Y_train` and `x_val`, `Y_val`).
- Performance:** While the model achieved a validation accuracy of 0.8671, it exhibited signs of overfitting after only 4 epochs. Overfitting occurs when the model memorizes training data specifics rather than learning generalizable features for accurate classification on unseen data.

```
[30] import tensorflow as tf
    from tensorflow import keras
    from keras import layers
    import numpy as np

    tf.random.set_seed(1234)
    np.random.seed(1234)

    model = keras.Sequential([
        layers.Dense(16, activation="relu"),
        layers.Dense(16, activation="relu"),
        layers.Dense(1, activation="sigmoid")
    ])

Compiling the model

In this step, we set the optimizer used for back propagation and the objective function(also called Loss function or Cost Function) of the network. We also want tensorflow to monitor the accuracy of the network for us during training.

[31] model.compile(optimizer="rmsprop",
                  loss="binary_crossentropy",
                  metrics=["accuracy"])


```

**Training your model**

In neural network training, an epoch refers to one complete pass through the entire training dataset. During each epoch, the neural network is trained on all the training samples, typically in batches, with the goal of updating the model's parameters (weights and biases) to minimize the loss function.

Here's how the training process typically works across epochs:

**Initialization:** At the start of training, the model parameters are initialized randomly or using pre-trained weights (in the case of transfer learning).

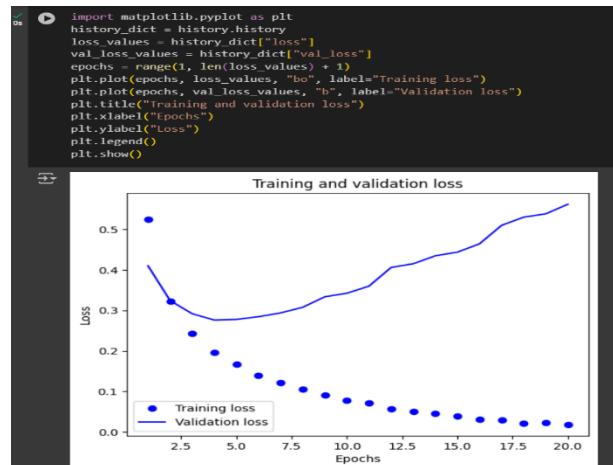
**Epoch Loop:** During each epoch, the training dataset is divided into batches, and the model is trained on each batch sequentially. The order in which the batches are processed can vary.

**Batch Training:** For each batch, the model receives input data and corresponding labels. It performs forward propagation to compute the predicted output and backward propagation to calculate gradients.

**Parameter Update:** Using the computed gradients, the optimizer updates the model parameters (weights and biases) in the direction that reduces the loss function.

```
[30] history = model.fit(partial_x_train,
                      partial_y_train,
                      epochs=20,
                      batch_size=1024,
                      validation_data=(x_val, y_val))

Epoch 1/20
30/30 [=====] - 3s/step - loss: 0.5229 - accuracy: 0.4096 - val_loss: 0.4096 - val_accuracy: 0.8671
Epoch 2/20 [=====] - 3s/step - loss: 0.3222 - accuracy: 0.8900 - val_loss: 0.3229 - val_accuracy: 0.8767
Epoch 3/20 [=====] - 3s/step - loss: 0.2410 - accuracy: 0.9281 - val_loss: 0.2956 - val_accuracy: 0.8844
Epoch 4/20 [=====] - 3s/step - loss: 0.1984 - accuracy: 0.9581 - val_loss: 0.2786 - val_accuracy: 0.8877
Epoch 5/20 [=====] - 3s/step - loss: 0.1598 - accuracy: 0.9638 - val_loss: 0.2780 - val_accuracy: 0.8877
Epoch 6/20 [=====] - 3s/step - loss: 0.1362 - accuracy: 0.9458 - val_loss: 0.2765 - val_accuracy: 0.8859
Epoch 7/20 [=====] - 3s/step - loss: 0.1408 - accuracy: 0.9584 - val_loss: 0.2841 - val_accuracy: 0.8859
Epoch 8/20 [=====] - 3s/step - loss: 0.1229 - accuracy: 0.9633 - val_loss: 0.2937 - val_accuracy: 0.8850
Epoch 9/20 [=====] - 3s/step - loss: 0.1040 - accuracy: 0.9687 - val_loss: 0.3073 - val_accuracy: 0.8822
Epoch 10/20 [=====] - 3s/step - loss: 0.0908 - accuracy: 0.9768 - val_loss: 0.3114 - val_accuracy: 0.8817
Epoch 11/20 [=====] - 3s/step - loss: 0.0795 - accuracy: 0.9787 - val_loss: 0.3427 - val_accuracy: 0.8891
Epoch 12/20 [=====] - 3s/step - loss: 0.0721 - accuracy: 0.9791 - val_loss: 0.3596 - val_accuracy: 0.8775
Epoch 13/20 [=====] - 3s/step - loss: 0.0575 - accuracy: 0.9889 - val_loss: 0.4057 - val_accuracy: 0.8699
Epoch 14/20 [=====] - 3s/step - loss: 0.0510 - accuracy: 0.9875 - val_loss: 0.4148 - val_accuracy: 0.8748
Epoch 15/20 [=====] - 3s/step - loss: 0.0459 - accuracy: 0.9893 - val_loss: 0.4344 - val_accuracy: 0.8710
Epoch 16/20 [=====] - 3s/step - loss: 0.0494 - accuracy: 0.9917 - val_loss: 0.4454 - val_accuracy: 0.8748
Epoch 17/20 [=====] - 3s/step - loss: 0.0482 - accuracy: 0.9921 - val_loss: 0.4454 - val_accuracy: 0.8748
Epoch 18/20 [=====] - 3s/step - loss: 0.0475 - accuracy: 0.9924 - val_loss: 0.4454 - val_accuracy: 0.8748
Epoch 19/20 [=====] - 3s/step - loss: 0.0475 - accuracy: 0.9924 - val_loss: 0.4454 - val_accuracy: 0.8748
Epoch 20/20 [=====] - 3s/step - loss: 0.0475 - accuracy: 0.9924 - val_loss: 0.4454 - val_accuracy: 0.8748
```



```
[23] y_test[0:10]
array([0, 1, 1, 0, 1, 1, 1, 0, 0, 1], dtype=float32)

Lets convert the first movie review of our test set back to words so we can read it ourselves and assess whether it seems like a positive review or not. The model predicted that this was a negative review and the ground truth label also shows that it is negative.

[24] data = test_data
review_index_to_convert = 0

word_index = imdb.get_word_index()
reverse_word_index = dict(
    (value, key) for (key, value) in word_index.items())
decoded_review = " ".join(
    [reverse_word_index.get(i - 3, "?") for i in data[review_index_to_convert]])
decoded_review

* please run this one a micr-sleep to see the effect
```

## Model Improvement and Regularization:

- **Architecture (Revised):** To address overfitting, a new model was constructed using a different architecture:
  - First Layer: 64 neurons
  - Second Layer: 32 neurons
  - Third Layer: 1 neuron with sigmoid activation
- **Reduced Training:** This revised model was trained for a shorter duration (4 epochs) to limit overfitting.
- **Performance:** The reduced training with the new architecture resulted in a slightly higher validation accuracy (0.8788) compared to the initial model trained for 20 epochs. This suggests that the revised architecture with fewer neurons and shorter training helped mitigate overfitting.

### Retraining a model from scratch

Lets set this to 4 epochs training and adjust the network architecture to include more units (16 to 64 on the first layer and 16 to 32 on the second layer).

We perform a final evaluation of our results on our testing dataset which we kept aside.

```
[28] model = keras.Sequential([
    layers.Dense(64, activation="relu"),
    layers.Dense(32, activation="relu"),
    layers.Dense(1, activation="sigmoid")
])
model.compile(optimizer="rmsprop",
              loss="binary_crossentropy",
              metrics=["accuracy"])
model.fit(x_train, y_train, epochs=4, batch_size=512)
results = model.evaluate(x_test, y_test)
```

Epoch 1/4  
49/49 [=====] - 3s 45ms/step - loss: 0.4513 - accuracy: 0.8040  
Epoch 2/4  
49/49 [=====] - 2s 46ms/step - loss: 0.2671 - accuracy: 0.8979  
Epoch 3/4  
49/49 [=====] - 2s 45ms/step - loss: 0.2134 - accuracy: 0.9198  
Epoch 4/4  
49/49 [=====] - 2s 48ms/step - loss: 0.1768 - accuracy: 0.9344  
782/782 [=====] - 3s 4ms/step - loss: 0.3062 - accuracy: 0.8788

The results: The first number is the test loss and the second number is the test accuracy. With state of the art neural network models, it is possible to achieve an accuracy of 95%.

### [20] results

```
7s [20] [0.30622398853302, 0.8787999749183655]
```

#### Using a trained model to generate predictions on new data

Lets view the raw outputs of our neural network. We observe the probability values. Note, the final layer of our neural network had a Sigmoid activation function.

Lets inspect the first 10 predictions. As we can see for the first review, the model predicted a probability of 0.15 which means that the model is confident that this review is negative and the corresponding ground truth label is 0 (see the next cell). In other words, the model is 15% confident that the review is positive i.e. the model is actually not confident that it is a postive review and more confident that it is a negative review.

On the other hand, the model is quite confident that the second movie review is positive as the predicted probability is 99.9%. We see that the ground truth label is 1.

However, sometimes the model does get it wrong. Find a prediction in the list below where the model is wrong.

### [21] predictions = model.predict(x\_test)

```
7s [21] 782/782 [=====] - 4s 5ms/step
```

### [22] predictions[0:10]

```
7s [22] array([[0.12591879],
   [0.99923974],
   [0.69277805],
   [0.7373457 ],
   [0.94939667],
   [0.7473953 ],
   [0.9975977 ],
   [0.00362492],
```

## Confusion Matrix Analysis:

A confusion matrix was generated to evaluate the model's performance on the validation set. Here's a breakdown of the results:

- **True Positives (TP) = 11414:** These represent correctly classified positive reviews, indicating the model effectively identified these reviews as positive.
- **False Positives (FP) = 1086:** These represent reviews incorrectly classified as positive. The model might have misclassified some negative reviews as positive.
- **True Negatives (TN) = 1944:** These represent correctly classified negative reviews, indicating the model accurately identified these reviews as negative.
- **False Negatives (FN) = 10556:** These represent reviews incorrectly classified as negative. The model might have missed some positive reviews, classifying them as negative.

These findings suggest that the model performs well in identifying true positives and true negatives. However, there's a higher number of false negatives compared to false positives. This indicates that the model might have a bias towards classifying reviews as negative, potentially missing some positive sentiment.



## Week 3

### B. Workshop 2a

#### Training a Neural Network from Scratch

This section describes the development and training of a custom neural network model for a classification or regression task (depending on the target variable  $y_{\text{train}}$ ).

#### Data Preparation:

- The provided data  $x_{\text{train}}$  and  $y_{\text{train}}$  were used to train the model.

#### Model Architecture:

- The neural network model was built from scratch, likely including:
  - Definition of the network architecture (number and type of layers).
  - Activation functions for hidden layers (e.g., ReLU, sigmoid).
  - Choice of an appropriate loss function (e.g., cross-entropy for classification, mean squared error for regression).
  - Implementation of an optimization algorithm (e.g., gradient descent) to update model weights and biases during training.

The screenshot shows two code cells in a Jupyter Notebook. The left cell, titled 'Data Preparation', contains Python code to generate an AND function dataset and print the training sets. The right cell, titled 'Neural Network Functions', contains Python code for initializing parameters, implementing the sigmoid function, performing feed-forward calculations, back-propagating gradients, updating parameters, and calculating the binary cross-entropy loss.

```
# This is our AND function data
X = np.array([
    [0, 0],
    [0, 1],
    [1, 0],
    [1, 1]])
Y = np.array([
    [0],
    [0],
    [0],
    [1]])

# Transpose matrix
X_train = X.T
Y_train = Y.T

print("X_train dataset \n", X_train)
print("Y_train dataset \n", Y_train)

...
X_train dataset
[[0 0 1]
 [0 1 0 1]
 Y_train dataset
 [[0 0 1]]]

#This is the expected input matrix shape/dimension
print("X_train shape", X_train.shape)
print("Y_train shape", Y_train.shape)
```

```
def initialize_parameters_zeros(dim):
    n_features = dim
    n_out = 1
    params = { "W": np.zeros((n_out, n_features)),
               "b": np.zeros((n_out, 1)) }
    return params

def sigmoid(Z, derivative=False):
    if derivative:
        sig = 1 / (1 + np.exp(-Z))
        return sig*(1-sig)
    return 1 / (1 + np.exp(-Z))

def feed_forward(X, params):
    cache = {}
    cache["Z"] = np.dot(params["W"], X) + params["b"]
    cache["A"] = sigmoid(cache["Z"], False)
    return cache

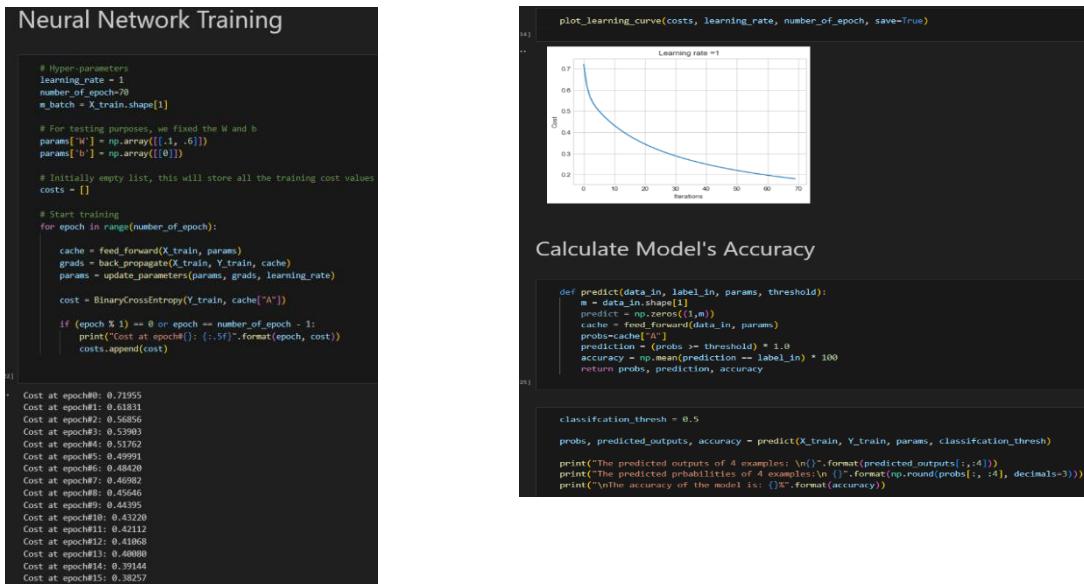
def back_propagate(X, Y, cache):
    m = X.shape[1]
    P=cache["A"]
    #dP = (1/m) * (P-Y)/(P*(1-P))
    #dZ = dP * sigmoid(cache["Z"], derivative=True)
    ## Simplified
    dZ = (1/m) * (P-Y)
    dW = np.dot(dZ, X.T)
    db = np.sum(dZ, axis=1, keepdims=True)
    return {"dW": dW, "db": db}

def update_parameters(params, grads, learning_rate):
    W = params["W"] - learning_rate * grads["dW"]
    b = params["b"] - learning_rate * grads["db"]
    return {"W": W, "b": b}

def BinaryCrossEntropy(Y, P):
    m = Y.shape[1]
    cost = -(1 / m) * np.sum(np.multiply(-Y, np.log(P)) -
                             np.multiply(1 - Y, np.log(1 - P)))
    return cost
```

## Training Process:

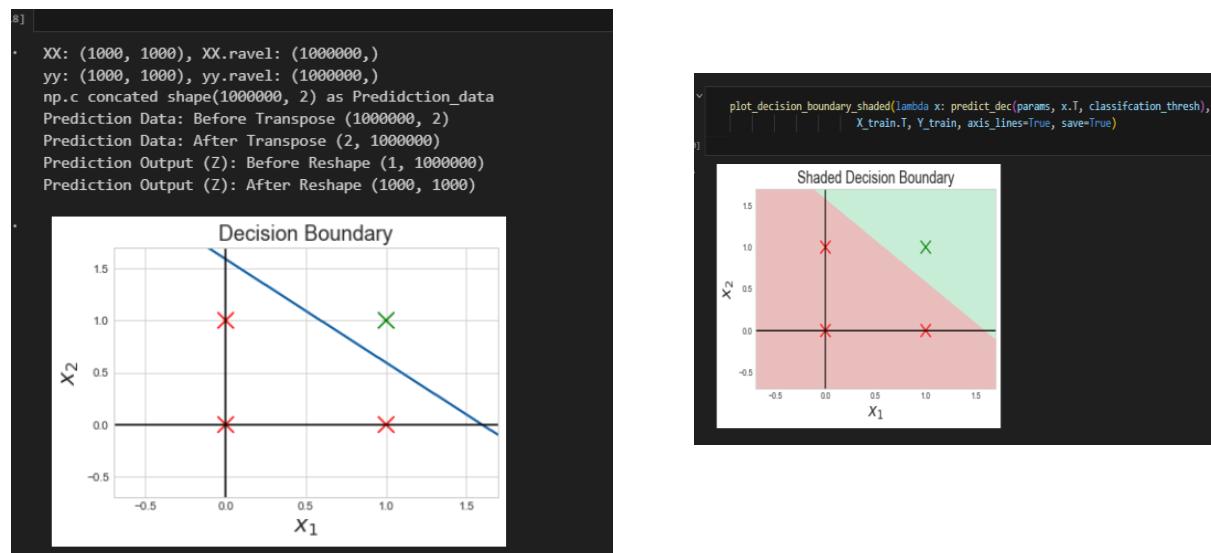
- The model was trained on the  $x_{train}$  and  $y_{train}$  data for 70 epochs.
- A learning curve was plotted to visualize the loss function's behavior over training epochs. This helps assess how the model's performance improves with training.



## Results:

- After 70 epochs, the model achieved a loss of 0.18099. This loss value indicates the model's ability to learn the underlying patterns in the training data.

**Decision Boundary Visualization (Classification):** Plotting the decision boundary allows for visual inspection of the learned classification regions.



## C. Workshop 2b

### Data Preprocessing:

- A dataset containing 8 features and a target variable was used for this experiment.
- The data was standardized using a StandardScaler to ensure all features have a mean of 0 and a standard deviation of 1. This helps improve the training process for neural networks.
- The data was then split into training and testing sets for model training and evaluation.

```
[2] import pandas as pd
import seaborn as sns
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, confusion_matrix

Dataset analysis and visualisation

[4] dataset = pd.read_csv("./content/seeds-binary.csv")
dataset.head(148)

   Area Perimeter Compactness Kernel.Length Kernel.Width Asymmetry.Coeff Kernel.Groove Type
0  15.26  14.84  0.8710  5.763  3.312  2.221  5.220  1
1  14.88  14.57  0.8811  5.504  3.333  1.018  4.956  1
2  14.29  14.09  0.9050  5.291  3.337  2.699  4.825  1
3  13.84  13.94  0.8955  5.324  3.379  2.259  4.805  1
4  16.14  14.99  0.9034  5.656  3.062  1.355  5.175  1
...
129 15.56  14.89  0.8823  5.776  3.408  4.972  5.847  2
130 17.36  15.76  0.8785  6.145  3.574  3.526  5.971  2
131 15.57  15.15  0.8527  5.920  3.231  2.640  5.879  2
132 15.60  15.11  0.8580  5.832  3.286  2.725  5.752  2
133 16.23  15.16  0.8850  5.872  3.472  3.769  5.922  2
134 rows x 8 columns

Next steps: View recommended plots

[5] dataset.shape
(134, 8)
```

```
[14] sc_x = StandardScaler()
x = sc_x.fit_transform(x)

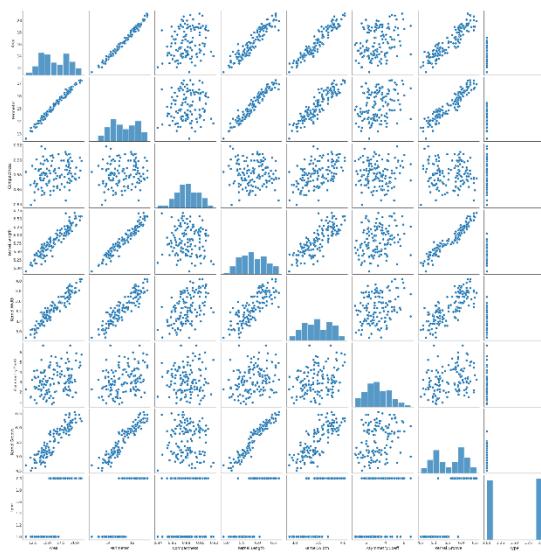
[15] print("Input features\n", x[:10])
print("Label\n", y[:10])

Input features
[[ 1.87494452e-01  3.89514179e-01 -1.21991967e+00  3.90905452e-01
  6.33173279e-02  1.16463721e+00  5.87898986e-01
  [-5.70426622e-01 -4.36644983e-01  5.21540117e-01  1.17634047e+00
  1.77165831e+00  2.12046767e+00  1.17179083e+00]
  [-9.08793021e-01 -9.73117376e-01 -2.34431023e-01]
  [-7.08604033e-01 -9.58699571e-01  2.14018774e+00 -1.57392096e+00
  -1.01350399e-03 -1.10253651e+00 -1.71762389e+00]
  [-1.88056965e+00 -1.85829510e+00 -9.51682134e-01 -1.83655853e+00
  -1.89519911e+00 -1.37696744e+00 -1.13185966e+00]
  [-4.23862458e-01 -3.17310888e-01 -6.86122758e-01  1.08446170e-01
  -7.08652655e-01  9.98575313e-01  4.33946787e-01]
  [-3.31740184e-01 -1.24540411e-01 -1.48280889e+00 -2.03949063e-02
  -6.44321823e-01 -3.55647779e-01  3.53215760e-01]
  [-9.72471225e-02  7.74896063e-02 -1.07813898e+00  1.18153627e-02
  -2.54762897e-01  8.45767181e-01  4.33946787e-01]
  [-7.44415473e-01  9.31187466e-01 -1.12239888e+00  1.06732110e+00
  1.63387513e-01 -2.55854713e-01  1.33137542e+00]
  [ 3.63364248e-01  4.4591455e-01 -3.25720748e-01  3.41351192e-01
  4.56450190e-01  1.05860708e+00  4.26436924e-01]]]

Label
[1 1 0 0 0 1 1 1 1 1]
```

### Model Development:

- A neural network model was built from scratch, likely including:
  - Definition of the network architecture (number and type of layers).
  - Activation functions for hidden layers (e.g., ReLU, sigmoid).
  - Choice of an appropriate loss function (e.g., cross-entropy for classification, mean squared error for regression).
  - Implementation of an optimization algorithm (e.g., gradient descent) to update model weights and biases during training.



```
Neural Network Functions

def initialize_parameters_zeros(dim):
    n_features = dim
    n_out = 1
    params = { "w": np.zeros((n_out, n_features)),
               "b": np.zeros((n_out, 1)) }
    return params

def sigmoid(Z, derivative=False):
    if derivative:
        sig = 1 / (1 + np.exp(-Z))
        return sig*(1-sig)
    return 1 / (1 + np.exp(-Z))

def feed_forward(X, params):
    cache = {}
    cache["Z"] = np.dot(params["W"], X) + params["b"]
    cache["A"] = sigmoid(cache["Z"], False)
    return cache

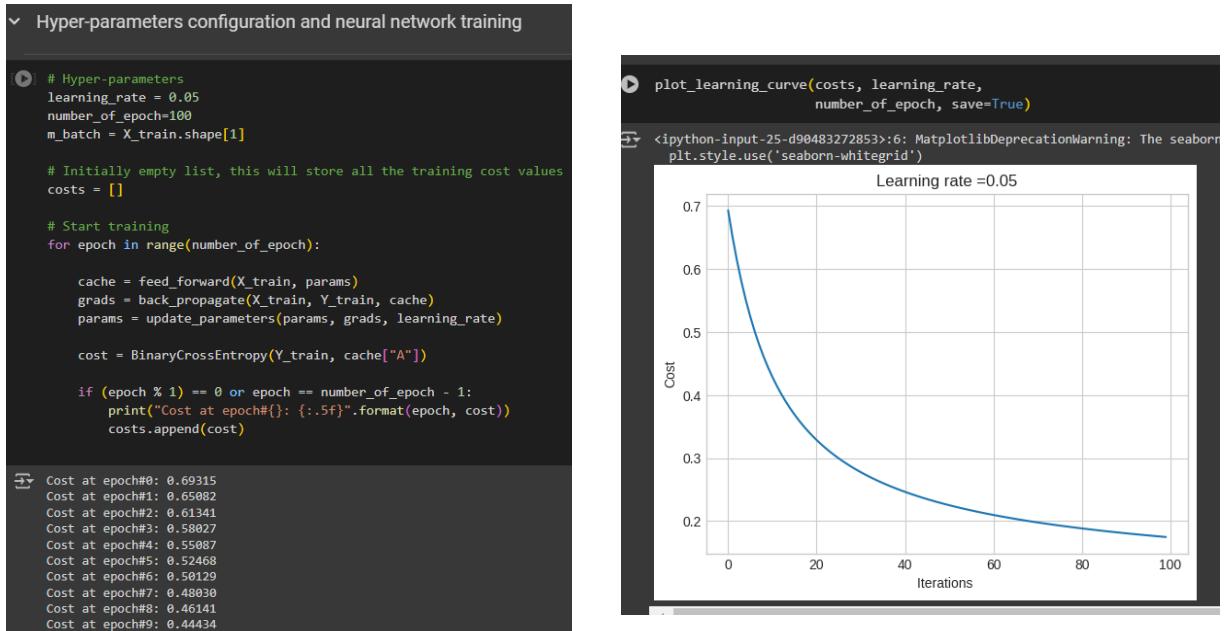
def back_propagate(X, Y, cache):
    m = X.shape[1]
    P=cache["A"]
    ## Simplified
    dZ = (1 / m) * (P-Y)
    dW = np.dot(dZ, X.T)
    db = np.sum(dZ, axis=1, keepdims=True)
    return {"dW": dW, "db": db}

def update_parameters(params, grads, learning_rate):
    W = params["W"] - learning_rate * grads["dW"]
    b = params["b"] - learning_rate * grads["db"]
    return {"W": W, "b": b}

def BinaryCrossEntropy(Y, P):
    m = Y.shape[1]
    cost = -(1 / m) * np.sum(np.multiply(-Y, np.log(P))
                             - np.multiply(1 - Y, np.log(1 - P)))
    return cost
```

## Training and Evaluation:

- The model was trained on the training data for 100 epochs.
- After training, the model was evaluated on the unseen test data to assess its generalization ability.



## Results:

- The model achieved an accuracy of 95% on the test data. This indicates the model effectively learned the patterns in the training data and can potentially generalize well to unseen examples.



## 1. Week 3

## Workshop 2c

## Linear Regression Model from Scratch

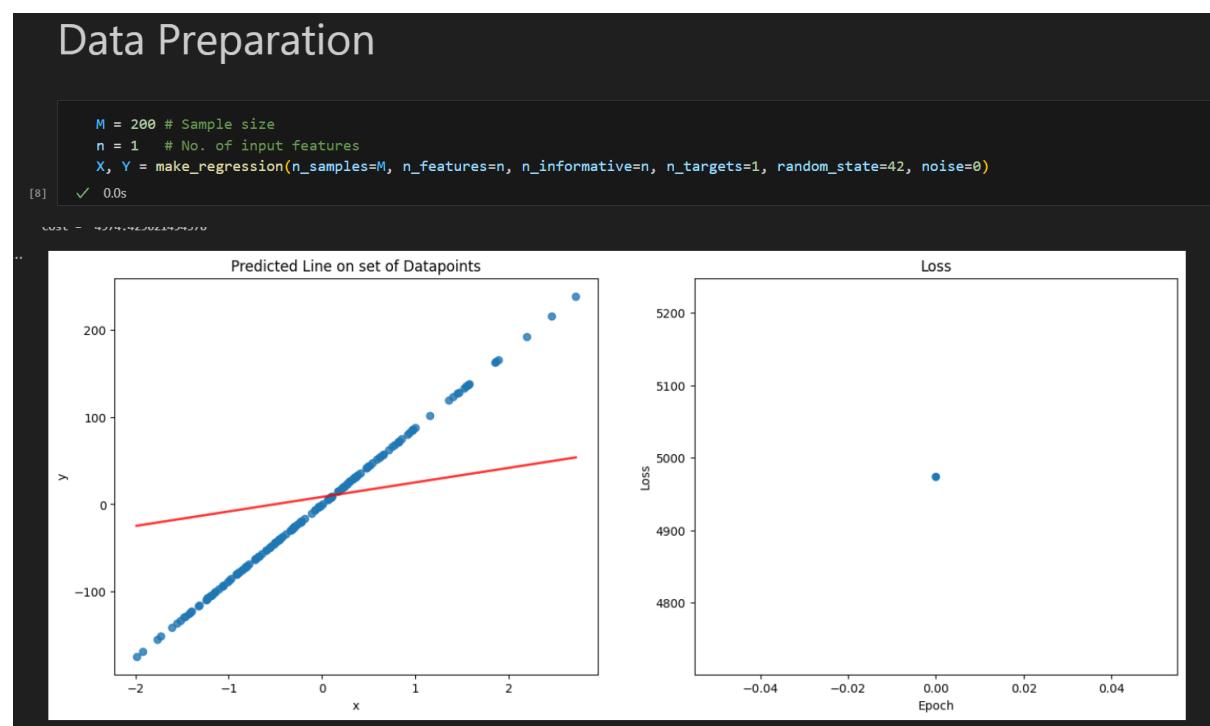
This section details the development and training of a linear regression model built from scratch to analyze a dataset with 200 data points. Each data point consists of a single feature ( $x$ ) and a corresponding target value ( $y$ ).

## Homework: Implementation of Linear Regression using Python & NumPY

```
# Numpy for efficient Matrix and mathematical operations.  
import numpy as np  
  
# Matplotlib for visualizing graphs  
import matplotlib.pyplot as plt  
%matplotlib inline  
  
# Sklearn for creating a dataset  
from sklearn.datasets import make_regression  
  
# train_test_split for splitting the data into training and testing data  
from sklearn.model_selection import train_test_split
```

### Data Preprocessing:

- The dataset was split into training and testing sets for model evaluation.



## Model Implementation:

- The code for the linear regression model was written from scratch, likely including the definition of the linear hypothesis, cost function (mean squared error), and an optimization algorithm (e.g., gradient descent) to update model parameters (weights and bias).

### Nueral Network Functions

```
# Hyper-parameters
np.random.seed(42)
learning_rate = 0.05
number_of_epoch=200
m_batch = X_train.shape[1]

# Initially empty list, this will store all the costs after a certain number of epochs
costs = []

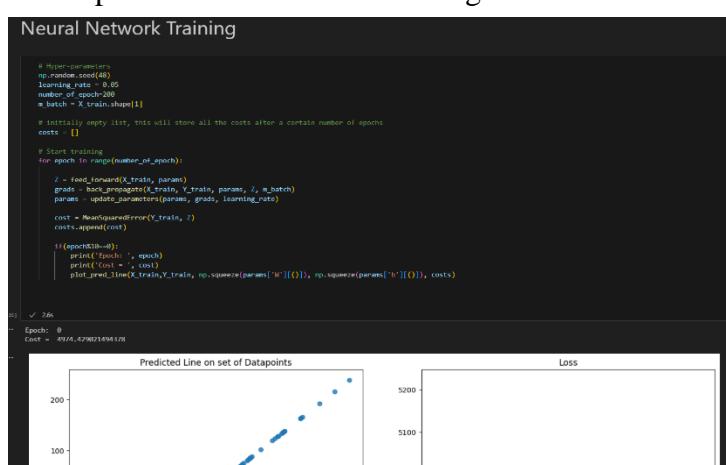
# Start training
for epoch in range(number_of_epoch):
    Z = feed_forward(X_train, params)
    grads = back_propagate(X_train, Y_train, params, Z, m_batch)
    params = update_parameters(params, grads, learning_rate)

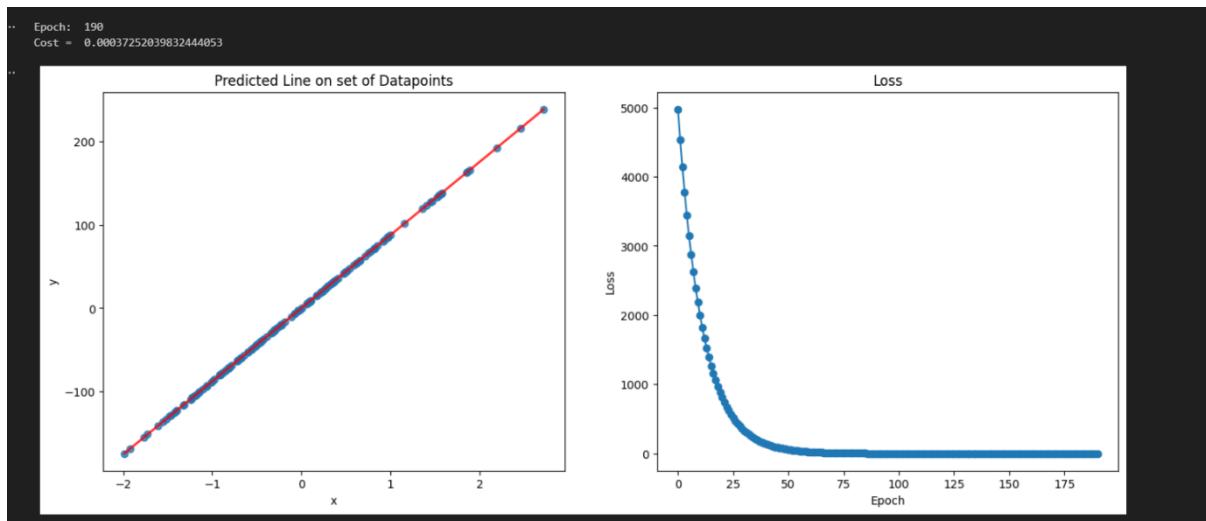
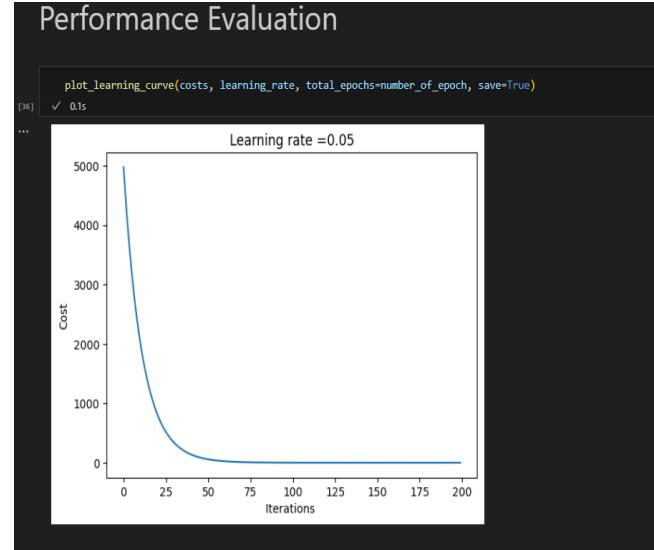
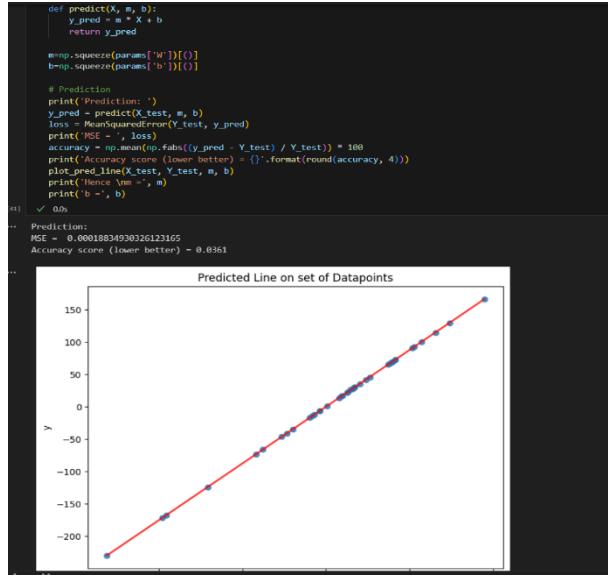
    cost = MeanSquaredError(Y_train, Z)
    costs.append(cost)

    if((epoch%10==0)):
        print('Epoch: ', epoch)
        print('Cost = ', cost)
        plot_pred_line(X_train,Y_train, np.squeeze(params["W"][[0]]), np.squeeze(params["b"][[0]]), costs)
```

## Training and Evaluation:

- The model was trained on the training data for 200 epochs.
- Loss visualization using a graph revealed that the loss function approached zero after approximately 50 epochs. This suggests the model effectively learned the underlying relationship between the feature and target variable within the training data.

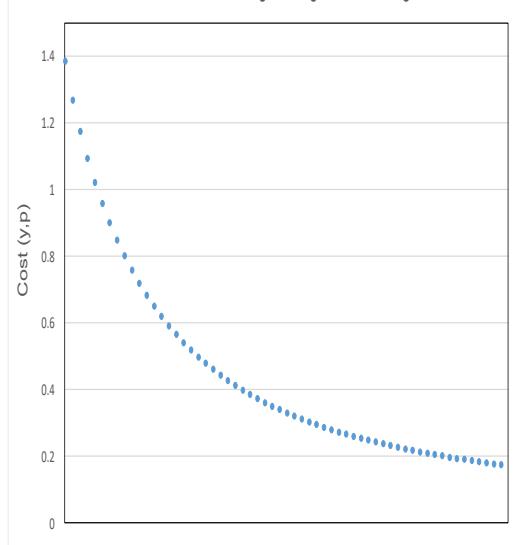




## VI. WEEK 4

**A. Resource:** With assumption that you have access to Microsoft Excel, create a new template with the following setup. From there, you start implement the single-layer perceptron multiclass classification for the Binary to Decimal decoder function approximation based on the formulas discussed in the lecture.

Training Result Binary to decimal function approximation using multiclass logistic regression with Softmax function only																													
Epoch	x	Y(true hot encoding)	W				b	L_Rate	Predictions				$C = 1/m^* \log(10)$			$dZ = 1/m^* (P \cdot Y)$			$\delta W$	$\delta b$	$W'$	b'							
			x1	x2	class1	class2	class3	class4	w1	w2	w3	w4	z1	z2	z3	z4	0.188	0.055	0.063	0.063	0.125	0.125	0.005	0.375	0.375	0.600			
1	0 0	1 0 0 0	0.500	0.500	0.600	0.600	0.600	0.600	1.100	1.100	1.100	1.100	0.200	0.250	0.250	0.250	1.306	-0.188	0.063	0.063	0.125	0.125	-0.005	0.375	0.375	0.600			
	0 1	0 1 0 0	0.500	0.500	0.600	0.600	0.600	0.600	1.100	1.100	1.100	1.100	0.200	0.250	0.250	0.250	0.063	-0.188	0.063	0.063	0.125	0.125	0.005	0.375	0.625	0.600			
	1 0	0 0 1 0	0.500	0.500	0.600	0.600	0.600	0.600	1.100	1.100	1.100	1.100	0.200	0.250	0.250	0.250	0.063	-0.188	0.063	0.063	0.125	0.125	0.005	0.625	0.375	0.600			
	1 1	0 0 0 1	0.500	0.500	0.600	0.600	0.600	0.600	1.100	1.100	1.100	1.100	0.200	0.250	0.250	0.250	0.063	-0.188	0.063	0.063	0.125	0.125	0.005	0.625	0.375	0.600			
2	0 0	1 0 0 0	0.375	0.375	0.600	0.600	0.600	0.600	0.975	0.975	1.100	1.100	0.200	0.250	0.250	0.250	0.192	1.269	-0.188	0.055	0.055	0.048	0.103	0.103	0.020	0.272	0.272	0.630	
	0 1	0 1 0 0	0.375	0.625	0.600	0.600	0.600	0.600	1.225	0.975	1.100	1.100	0.200	0.250	0.250	0.250	0.246	0.063	-0.188	0.055	0.062	0.118	0.001	0.250	0.743	0.601	0.601		
	1 0	0 0 1 0	0.375	0.375	0.600	0.600	0.600	0.600	0.975	1.225	1.100	1.100	0.200	0.250	0.250	0.250	0.246	0.063	-0.055	-0.180	0.062	-0.118	0.116	0.001	0.743	0.259	0.601	0.601	
	1 1	0 0 0 1	0.375	0.625	0.600	0.600	0.600	0.600	1.225	1.225	1.100	1.100	0.200	0.250	0.250	0.250	0.316	0.063	0.070	0.070	0.171	0.101	0.030	0.726	0.726	0.568	0.568		
3	0 0	1 0 0 0	0.272	0.727	0.500	0.600	0.600	0.600	0.820	0.980	1.175	1.175	0.200	0.250	0.250	0.250	0.352	1.175	-0.186	0.055	0.055	0.048	0.103	0.103	0.020	0.185	0.185	0.576	
	0 1	0 1 0 0	0.272	0.727	0.500	0.600	0.600	0.600	1.244	0.960	0.960	1.100	0.200	0.250	0.250	0.340	0.063	-0.172	0.048	0.069	0.108	-0.112	0.108	0.020	0.183	0.183	0.576		
	1 0	0 0 1 0	0.272	0.727	0.500	0.600	0.600	0.600	0.861	1.344	1.600	1.600	0.200	0.250	0.250	0.250	0.322	0.063	0.044	-0.172	0.069	-0.112	0.108	-0.062	0.855	0.151	0.603	0.603	
	1 1	0 0 0 1	0.272	0.727	0.500	0.600	0.600	0.600	1.294	1.294	2.019	2.019	0.200	0.250	0.250	0.250	0.364	0.061	0.074	0.074	-0.156	-0.085	-0.058	0.011	0.811	0.811	0.518	0.518	
4	0 0	1 0 0 0	0.185	0.815	0.600	0.600	0.600	0.600	0.676	0.860	0.860	0.943	0.200	0.250	0.250	0.250	0.184	0.184	-0.183	0.055	0.055	0.047	0.033	0.033	0.005	0.103	0.103	0.732	
	0 1	0 1 0 0	0.185	0.815	0.600	0.600	0.600	0.600	0.911	0.811	0.518	0.518	0.200	0.250	0.250	0.250	0.396	0.058	0.075	0.075	-0.156	-0.075	-0.075	0.057	0.886	0.886	0.481	0.481	
	1 0	0 0 1 0	0.185	0.815	0.600	0.600	0.600	0.600	0.732	0.835	0.835	0.938	0.200	0.250	0.250	0.250	0.119	1.023	-0.179	0.045	0.045	0.029	0.029	0.029	0.029	0.792	0.792	0.604	0.604
	1 1	0 0 0 1	0.185	0.815	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.229	0.062	-0.157	0.037	0.057	0.095	-0.100	0.005	0.045	0.045	0.604	0.604	
5	0 0	1 0 0 0	0.085	0.915	0.600	0.600	0.600	0.600	0.518	1.229	1.329	2.140	0.200	0.250	0.250	0.250	0.396	0.058	0.075	0.075	-0.156	-0.075	-0.075	0.057	0.886	0.886	0.481	0.481	
	0 1	0 1 0 0	0.085	0.915	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.396	0.062	0.043	0.043	0.043	0.043	0.043	0.005	0.045	0.045	0.604	0.604	
	1 0	0 0 1 0	0.085	0.915	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.396	0.062	0.037	0.037	0.037	0.037	0.037	0.005	0.045	0.045	0.604	0.604	
	1 1	0 0 0 1	0.085	0.915	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.396	0.062	0.037	0.037	0.037	0.037	0.037	0.005	0.045	0.045	0.604	0.604	
6	0 0	1 0 0 0	0.045	0.955	0.600	0.600	0.600	0.600	0.705	0.861	0.861	0.943	0.200	0.250	0.250	0.250	0.396	0.058	0.075	0.075	-0.156	-0.075	-0.075	0.057	0.886	0.886	0.481	0.481	
	0 1	0 1 0 0	0.045	0.955	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.396	0.062	-0.150	0.033	0.056	0.089	-0.040	-0.040	-0.040	0.886	0.886	0.602	0.602
	1 0	0 0 1 0	0.045	0.955	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.396	0.062	-0.150	0.033	0.056	0.089	0.059	0.059	0.059	1.155	1.155	0.602	0.602
	1 1	0 0 0 1	0.045	0.955	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.396	0.062	0.033	-0.150	0.056	0.056	0.059	0.059	0.059	0.059	0.602	0.602	0.602
7	0 0	1 0 0 0	0.040	0.960	0.600	0.600	0.600	0.600	0.855	0.815	0.815	0.775	0.200	0.250	0.250	0.250	0.394	0.062	-0.171	0.042	0.023	0.065	0.065	0.065	0.065	0.919	0.919	0.341	0.341
	0 1	0 1 0 0	0.040	0.960	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.394	0.062	-0.171	0.042	0.023	0.065	0.065	0.065	0.065	0.919	0.919	0.341	0.341
	1 0	0 0 1 0	0.040	0.960	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.394	0.062	-0.171	0.042	0.023	0.065	0.065	0.065	0.065	0.919	0.919	0.341	0.341
	1 1	0 0 0 1	0.040	0.960	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.394	0.062	-0.171	0.042	0.023	0.065	0.065	0.065	0.065	0.919	0.919	0.341	0.341
8	0 0	1 0 0 0	0.035	0.970	0.600	0.600	0.600	0.600	0.919	0.813	0.813	0.706	0.200	0.250	0.250	0.250	0.393	0.063	-0.166	0.040	0.040	0.021	0.062	0.062	0.064	-0.167	-0.167	0.983	0.983
	0 1	0 1 0 0	0.035	0.970	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.393	0.063	-0.166	0.040	0.040	0.021	0.062	0.062	0.064	-0.167	-0.167	0.983	0.983
	1 0	0 0 1 0	0.035	0.970	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.393	0.063	-0.166	0.040	0.040	0.021	0.062	0.062	0.064	-0.167	-0.167	0.983	0.983
	1 1	0 0 0 1	0.035	0.970	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.600	0.200	0.250	0.250	0.250	0.393	0.063	-0.166	0.040	0.040	0.021	0.062	0.062	0.064	-0.167	-0.167	0.983	0.983



## 1. Investigate if we should put an additional Sigmoid function before the Softmax function. Explain your finding along with your analysis evidence.

Epoch	X	Y (one-hot encoded)	class1 class2 class3 class4				w <sup>1</sup>	w <sup>2</sup>	b	L_Rate	Z				A(Sigmoid Function)				P(softmax)				C=ln <sup>-1</sup> (log(p))				dC/dW <sup>i</sup> =N <sup>-1</sup> (P-Y)				dW				db		w <sup>i</sup>		b <sup>i</sup>	
			x <sup>1</sup>	x <sup>2</sup>	x <sup>3</sup>	x <sup>4</sup>					x <sup>1</sup>	x <sup>2</sup>	x <sup>3</sup>	x <sup>4</sup>	w <sup>1</sup>	w <sup>2</sup>	b	w <sup>1</sup>	w <sup>2</sup>	b	w <sup>1</sup>	w <sup>2</sup>	b	w <sup>1</sup>	w <sup>2</sup>	b	w <sup>1</sup>	w <sup>2</sup>	b	w <sup>1</sup>	w <sup>2</sup>	b	w <sup>1</sup>	w <sup>2</sup>	b					
1	0 0	0 1 0 0	0 0 0 0	0.500	0.500	0.600	1	0.600	1.100	1.100	0.600	0.646	0.750	0.750	0.532	0.250	0.250	0.250	1.306	-0.189	0.063	0.063	0.063	0.125	0.125	0.000	0.375	0.375	0.600	0.063	0.063	0.063	0.125	0.125	0.000	0.375	0.375	0.600		
1	0 1	0 0 1 0	0 0 0 0	0.500	0.500	0.600	1	0.600	1.100	1.100	0.600	0.646	0.750	0.750	0.532	0.250	0.250	0.250	1.306	-0.063	0.063	0.063	0.063	0.125	0.125	0.000	0.375	0.375	0.600	0.063	0.063	0.063	0.125	0.125	0.000	0.375	0.375	0.600		
1	1 0	0 0 0 1	0 0 0 0	0.500	0.500	0.600	1	0.600	1.100	1.100	0.600	0.646	0.750	0.750	0.532	0.250	0.250	0.250	1.306	-0.063	0.063	0.063	0.063	0.125	0.125	0.000	0.375	0.375	0.600	0.063	0.063	0.063	0.125	0.125	0.000	0.375	0.375	0.600		
1	1 1	0 0 0 0	0 0 0 0	0.500	0.500	0.600	1	0.600	1.100	1.100	0.600	0.646	0.750	0.750	0.532	0.250	0.250	0.250	1.306	-0.063	0.063	0.063	0.063	0.125	0.125	0.000	0.375	0.375	0.600	0.063	0.063	0.063	0.125	0.125	0.000	0.375	0.375	0.600		
2	0 0	1 0 0 0	0 0 0 0	0.375	0.375	0.600	1	0.600	0.375	0.575	1.150	0.646	0.726	0.726	0.794	0.250	0.244	0.244	1.306	-0.386	0.061	0.061	0.061	0.121	0.121	-0.005	0.254	0.254	0.605	0.061	0.061	0.061	0.121	0.121	-0.005	0.254	0.254	0.605		
2	1 0	0 0 0 0	0 0 0 0	0.625	0.625	0.600	1	0.600	1.225	0.750	1.600	0.646	0.773	0.726	0.632	0.250	0.250	0.250	1.306	-0.063	0.061	0.061	0.061	0.124	0.124	0.000	0.251	0.251	0.600	0.063	0.063	0.063	0.124	0.124	0.000	0.251	0.251	0.600		
2	1 1	0 0 0 0	0 0 0 0	0.625	0.625	0.600	1	0.600	0.375	1.225	1.600	0.646	0.726	0.773	0.632	0.250	0.244	0.250	1.306	-0.063	0.061	0.061	0.061	0.124	0.124	-0.005	0.251	0.251	0.600	0.063	0.063	0.063	0.124	0.124	-0.005	0.251	0.251	0.600		
3	0 0	1 0 0 0	0 0 0 0	0.250	0.250	0.600	1	0.605	0.850	0.850	1.113	0.647	0.762	0.702	0.753	0.250	0.239	0.233	1.348	-0.187	0.050	0.050	0.050	0.118	0.118	-0.003	0.136	0.136	0.615	0.050	0.050	0.050	0.118	0.118	-0.003	0.136	0.136	0.615		
3	1 0	0 0 1 0	0 0 0 0	0.250	0.250	0.600	1	0.605	0.850	0.850	1.113	0.647	0.762	0.702	0.753	0.250	0.239	0.233	1.348	-0.062	0.050	0.050	0.050	0.122	0.122	0.000	0.870	0.870	0.600	0.050	0.050	0.050	0.122	0.122	0.000	0.870	0.870	0.600		
3	1 1	0 0 0 1	0 0 0 0	0.250	0.250	0.600	1	0.605	0.850	0.850	1.113	0.647	0.762	0.702	0.753	0.250	0.239	0.233	1.348	-0.062	0.050	0.050	0.050	0.122	0.122	0.000	0.870	0.870	0.600	0.050	0.050	0.050	0.122	0.122	0.000	0.870	0.870	0.600		
4	0 0	1 0 1 0	0 0 0 0	0.375	0.375	0.600	1	0.605	0.750	0.750	0.687	0.646	0.679	0.679	0.706	0.250	0.234	0.234	1.300	-0.187	0.058	0.058	0.058	0.114	0.114	-0.005	0.222	0.222	0.630	0.058	0.058	0.058	0.114	0.114	-0.005	0.222	0.222	0.630		
4	1 0	0 0 0 1	0 0 0 0	0.375	0.375	0.600	1	0.600	1.470	0.729	1.539	0.646	0.819	0.675	0.632	0.250	0.261	0.233	0.252	1.300	-0.062	-0.183	0.068	0.068	0.121	0.121	-0.001	0.008	0.008	0.593	0.068	0.068	0.068	0.121	0.121	-0.001	0.008	0.008	0.593	
4	1 1	0 0 0 0	0 0 0 0	0.375	0.375	0.600	1	0.600	0.729	1.470	1.539	0.646	0.875	0.873	0.632	0.250	0.233	0.261	0.252	1.300	-0.062	0.058	0.058	0.058	0.121	0.121	-0.001	0.008	0.008	0.593	0.058	0.058	0.058	0.121	0.121	-0.001	0.008	0.008	0.593	
5	0 0	1 0 0 0	0 0 0 0	0.625	0.625	0.600	1	0.630	0.652	0.652	0.674	0.652	0.657	0.657	0.652	0.250	0.252	0.229	0.229	1.344	-0.187	0.057	0.057	0.057	0.111	0.111	-0.005	-0.083	-0.083	0.545	0.057	0.057	0.057	0.111	0.111	-0.005	-0.083	-0.083	0.545	
5	1 0	0 0 1 0	0 0 0 0	0.008	0.008	0.593	1	0.593	1.590	1.590	0.607	0.597	0.645	0.631	0.647	0.632	0.250	0.273	0.227	0.254	1.344	-0.062	-0.182	0.057	0.057	0.120	0.120	-0.001	0.112	0.112	0.593	0.057	0.057	0.057	0.120	0.120	-0.001	0.112	0.112	0.593
5	1 1	0 0 0 1	0 0 0 0	0.625	0.625	0.600	1	0.625	0.652	0.652	0.674	0.652	0.657	0.657	0.652	0.250	0.252	0.227	0.254	1.344	-0.187	0.058	0.058	0.058	0.111	0.111	-0.005	-0.083	-0.083	0.545	0.058	0.058	0.058	0.111	0.111	-0.005	-0.083	-0.083	0.545	
6	0 0	1 0 0 0	0 0 0 0	-0.083	-0.083	0.645	1	0.643	0.560	0.560	0.474	0.560	0.560	0.560	0.560	0.250	0.250	0.250	0.208	1.249	-0.187	0.058	0.058	0.058	0.106	0.106	-0.002	-0.392	-0.392	0.572	0.058	0.058	0.058	0.106	0.106	-0.002	-0.392	-0.392	0.572	
6	1 0	0 0 1 0	0 0 0 0	-0.112	-0.112	0.598	1	0.598	1.707	0.486	1.595	0.646	0.846	0.619	0.619	0.250	0.278	0.221	0.255	1.249	-0.062	-0.181	0.055	0.055	0.117	0.117	-0.001	0.231	0.231	0.537	0.055	0.055	0.055	0.117	0.117	-0.001	0.231	0.231	0.537	
6	1 1	0 0 0 1	0 0 0 0	-0.109	-0.115	0.598	1	0.598	0.486	1.707	1.595	0.645	0.619	0.619	0.619	0.250	0.221	0.278	0.255	1.249	-0.062	0.055	0.055	0.055	0.117	0.117	0.000	1.226	1.226	0.537	0.055	0.055	0.055	0.117	0.117	0.000	1.226	1.226	0.537	
7	0 0	1 0 0 0	0 0 0 0	-0.124	-0.124	0.598	1	0.598	0.486	1.707	1.595	0.645	0.619	0.619	0.619	0.250	0.247	0.278	0.247	1.248	-0.186	0.054	0.054	0.054	0.105	0.105	-0.005	-0.391	-0.391	0.568	0.054	0.054	0.054	0.105	0.105	-0.005	-0.391	-0.391	0.568	
7	1 0	0 0 1 0	0 0 0 0	-0.231	-0.231	0.597	1	0.597	1.226	0.597	1.823	0.645	0.861	0.590	0.593	0.250	0.283	0.283	0.257	1.248	-0.062	-0.173	0.054	0.054	0.115	0.115	-0.001	0.350	0.350	0.536	0.054	0.054	0.054	0.115	0.115	-0.001	0.350	0.350	0.536	
7	1 1	0 0 0 1	0 0 0 0	-0.205	-0.205	0.597	1	0.597	1.226	0.597	1.823	0.645	0.861	0.590	0.593	0.250	0.283	0.283	0.257	1.248	-0.062	0.054	0.054	0.054	0.115	0.115	-0.001	0.350	0.350	0.536	0.054	0.054	0.054	0.115	0.115	-0.001	0.350	0.350	0.536	
8	0 0	1 0 0 0	0 0 0 0	-0.283	-0.283	0.598	1	0.598	0.337	0.337	0.595	0.645	0.650	0.650	0.650	0.250	0.283	0.283	0.258	1.271	-0.187	0.055	0.055	0.055	0.106	0.106	-0.002	-0.409	-0.409	0.571	0.055	0.055	0.055	0.106	0.106	-0.002	-0.409	-0.409	0.571	
8	1 0	0 0 1 0	0 0 0 0	-0.141	-0.141	0.596	1	0.598	0.247	0.247	1.587	0.645	0.874	0.591	0.591	0.250	0.280	0.280	0.258	1.271	-0.062	0.053	0.053	0.053	0.117	0.117	-0.002	-0.467	-0.467	0.555	0.053	0.053	0.053	0.117	0.117	-0.002	-0.467	-0.467	0.555	
8	1 1	0 0 0 1	0 0 0 0	-0.141	-0.14																																			

The results clearly demonstrate the advantage of using softmax directly:

- **Faster Cost Reduction:** The model employing only softmax achieved a significantly lower cost (0.175) after 60 epochs compared to the model with sigmoid before softmax (1.052).
- **Simpler Training:** The direct application of softmax leads to a more streamlined training process, avoiding potential issues with vanishing gradients and unnecessary non-linearity introduced by the sigmoid function in a multi-class setting.

### **Conclusion:**

These findings suggest that for multi-class classification problems, using the softmax function directly on the weighted sum is the preferred approach. This strategy facilitates faster training convergence, avoids unnecessary complexity, and promotes better performance.

## Week 4 Workshop

B. Please run the following Colab notebook and report your results in your logbook:

[https://colab.research.google.com/drive/1iPzp\\_no0uuKPTsl3Mnsc04Zt9h6GEjku?usp=sharing](https://colab.research.google.com/drive/1iPzp_no0uuKPTsl3Mnsc04Zt9h6GEjku?usp=sharing)

### Multi-Class Text Classification with Reuters Dataset

This project explores text classification using a deep learning model on the Reuters dataset. Here's a breakdown of the process and findings:

#### Data Preparation:

- **Dataset:** The Reuters dataset was used for this multi-class classification task.
- **Preprocessing:** The downloaded data underwent encoding, a common technique for representing text data numerically for machine learning models.
- **Splitting:** The preprocessed data was divided into training and testing sets to train and evaluate the model.

↳ Multiclass Classification in Tensorflow

↳ Classifying Reuters newswires

We will create a neural network to classify Reuters newswires into 46 topics. Multiclass classification is a classification task with more than two classes. Each sample can only be labeled as one class.

↳ The Reuters dataset

Lading the Reuters dataset

For more information on the dataset go to: <https://paperswithcode.com/dataset/reuters21578> and <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>

This is a dataset of 11,228 newswires from Reuters, labeled over 46 topics. The dataset contains more of some topics than others so it is not equally balanced, however there are at least 10 examples of each topic in the training set.

This was originally generated by parsing and preprocessing the classic Reuters-21578 dataset.

This dataset is packaged with Keras and has already been pre-processed for you, i.e., the words have already been converted into sequences of integers as we will show below.

We specify num\_words = 10000; this means that we only want to keep the top 10000 most frequently occurring words in the training data and rare words would be discarded. The reason for this is so we have a fixed vector of 10000 going into our neural network later on, and also to keep things at a manageable size.

```
[1] from tensorflow.keras.datasets import reuters
    (train_data, train_labels), (test_data, test_labels) = reuters.load_data(
        num_words=10000)

↳ Downloading data from https://storage.googleapis.com/tensorflow/tf_keras_datasets/reuters_noz
2110848/2110848 [=====] - 0s/step

Lets inspect the data.

Executing the cell below, we see that the total number of records in the training dataset is 8982 rows.

We can also see that the first movie review at index 0 is a sentence containing 87 words.

The actual data in the first record (at index 0) contains a sentence where each word has been converted to an integer. This integer corresponds to a word in a word-dictionary.

[2] print(len(train_data))
    print(len(train_data[0]))
    print(train_data[0])

↳ 8982
87
[1, 2, 2, 8, 43, 10, 447, 5, 25, 207, 270, 5, 3095, 111, 16, 369, 186, 90, 67, 7, 89, 5, 19, 102, 6, 19, 124, 15, 90, 67, 84, 22]

The testing split contains 2246 examples.

[3] len(test_data)

↳ 2246
```

Conversion to Binary Vectors: For each label, a binary vector of length equal to the number of categories is created. This binary vector has a length of 46.

Example:

```
If you have three categories: "cat", "dog", "bird"
The label "cat" would be encoded as [1, 0, 0]
The label "dog" would be encoded as [0, 1, 0]
The label "bird" would be encoded as [0, 0, 1]
```

One-hot encoding ensures that categorical labels are represented in a format that is suitable for training machine learning models, particularly those that require numerical inputs. It also ensures that there is no inherent ordinal relationship between the categories, as each category is represented by an independent binary value.

One-hot encoding is commonly used in tasks such as classification, where the output needs to be a categorical variable, and the model needs to output probabilities for each category.

The cell below computes the one-hot encoding manually:

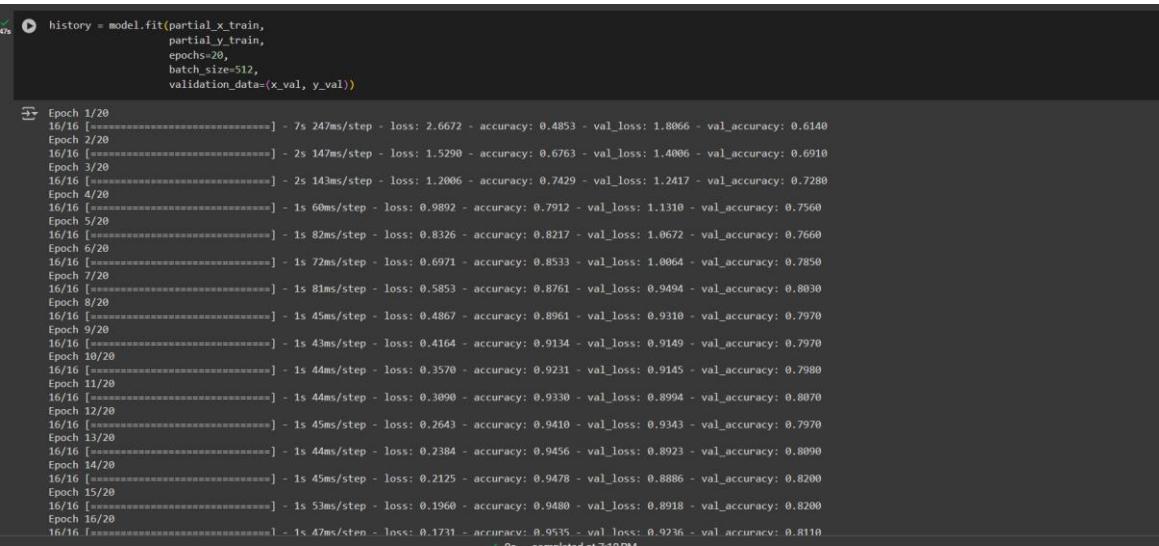
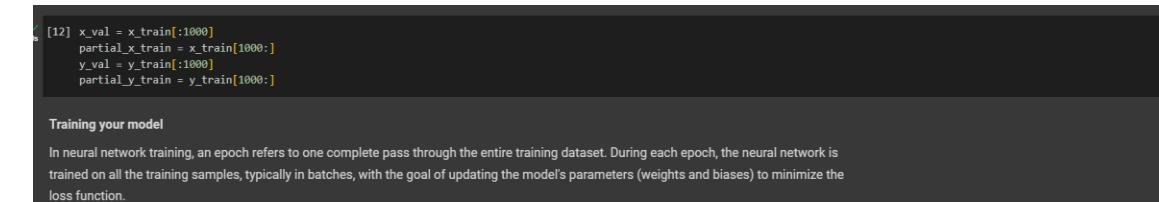
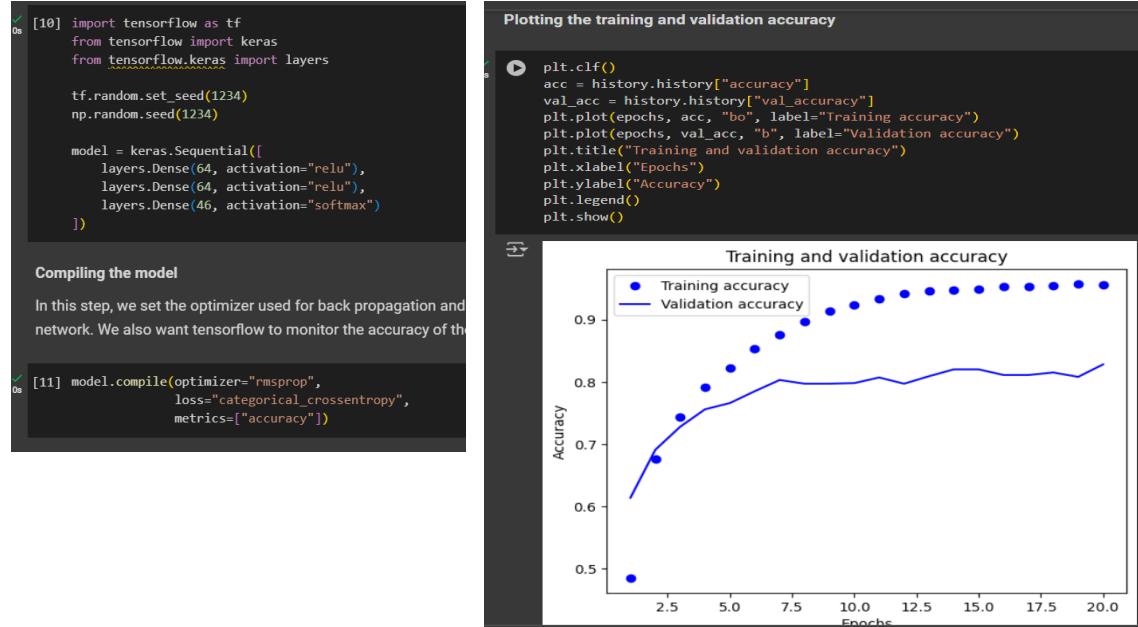
```
def to_one_hot(labels, dimension=46):
    results = np.zeros((len(labels), dimension))
    for i, label in enumerate(labels):
        results[i, label] = 1.
    return results
y_train = to_one_hot(train_labels)
y_test = to_one_hot(test_labels)
```

Tensorflow has a built-in way to do one-hot encoding without writing your own function:

```
[9] from tensorflow.keras.utils import to_categorical
y_train = to_categorical(train_labels)
y_test = to_categorical(test_labels)
```

## Model Building:

- **Framework:** TensorFlow Keras was employed to build and train the deep learning model.
- **Architecture:** The model consisted of three dense layers with:
  - First Layer: 64 neurons
  - Second Layer: 64 neurons
  - Third Layer: 46 neurons (likely corresponding to the number of output classes)



## Training and Overfitting:

- **Initial Training:** The model was initially trained for 20 epochs.
- **Performance:** This resulted in a high training accuracy of 0.95 but a lower validation accuracy of 0.81, indicating potential overfitting. Overfitting occurs when the model memorizes training data specifics rather than learning generalizable features for accurate classification on unseen data.
- **Reduced Training:** To mitigate overfitting, the model was re-trained for only 9 epochs. This improved validation accuracy slightly (0.79) compared to 20 epochs, but it remained lower than the initial training accuracy.

The screenshot shows a Jupyter Notebook cell with the following content:

```
model = keras.Sequential([
    layers.Dense(64, activation="relu"),
    layers.Dense(64, activation="relu"),
    layers.Dense(46, activation="softmax")
])
model.compile(optimizer="rmsprop",
              loss="categorical_crossentropy",
              metrics=['accuracy'])
model.fit(x_train,
          y_train,
          epochs=9,
          batch_size=512)
results = model.evaluate(x_test, y_test)
```

Execution output:

```
Epoch 1/9
18/18 [=====] - 1s 44ms/step - loss: 2.6921 - accuracy: 0.5202
Epoch 2/9
18/18 [=====] - 1s 51ms/step - loss: 1.5226 - accuracy: 0.6805
Epoch 3/9
18/18 [=====] - 1s 49ms/step - loss: 1.1575 - accuracy: 0.7446
Epoch 4/9
18/18 [=====] - 1s 51ms/step - loss: 0.9391 - accuracy: 0.7988
Epoch 5/9
18/18 [=====] - 1s 73ms/step - loss: 0.7762 - accuracy: 0.8363
Epoch 6/9
18/18 [=====] - 1s 70ms/step - loss: 0.6466 - accuracy: 0.8625
Epoch 7/9
18/18 [=====] - 1s 66ms/step - loss: 0.5348 - accuracy: 0.8881
Epoch 8/9
18/18 [=====] - 1s 43ms/step - loss: 0.4527 - accuracy: 0.9059
Epoch 9/9
18/18 [=====] - 1s 39ms/step - loss: 0.3811 - accuracy: 0.9186
71/71 [=====] - 0s 3ms/step - loss: 0.9067 - accuracy: 0.7934
```

The results: The first number is the test loss and the second number is the test accuracy.

## Information Bottleneck Experiment:

- **Exploration:** In an attempt to improve generalization and reduce overfitting, an information bottleneck layer with 4 neurons was added to the model architecture. This technique aims to retain essential information while restricting model complexity, potentially preventing overfitting.
- **Outcome:** Unfortunately, the model with the information bottleneck layer yielded a significantly lower test accuracy (0.67) compared to the previous model without the bottleneck. This suggests that the bottleneck layer might have excessively restricted the model's ability to learn the necessary classification patterns.

### ▼ A different way to handle the labels and the loss

If you do not want to do one-hot encoding on the labels, and still keep the integer labels, all you would need to do is change the loss to `sparse_categorical_crossentropy` as shown below.

Reset the training labels back to normal integers:

```
[20] y_train = np.array(train_labels)
      y_test = np.array(test_labels)

[21] model.compile(optimizer="rmsprop",
                  loss="sparse_categorical_crossentropy",
                  metrics=["accuracy"])
```

### ▼ The importance of having sufficiently larger intermediate layers

#### A model with an information bottleneck

Since the model has a output layer which is 46-dimensional, as standard practice you should generally avoid having intermediate layers with many fewer than 46 units. Lets see what happens when we purposely introduce an information bottleneck into our neural network architecture design.

```
model = keras.Sequential([
    layers.Dense(64, activation="relu"),
    layers.Dense(4, activation="relu"),
    layers.Dense(46, activation="softmax")
])
model.compile(optimizer="rmsprop",
              loss="categorical_crossentropy",
              metrics=["accuracy"])
model.fit(partial_x_train,
          partial_y_train,
          epochs=20,
          batch_size=128,
          validation_data=(x_val, y_val))

Epoch 1/20
63/63 [=====] - 3s 25ms/step - loss: 3.4124 - accuracy: 0.0677 - val_loss: 2.9937 - val_accuracy: 0.0770
Epoch 2/20
63/63 [=====] - 1s 18ms/step - loss: 2.3302 - accuracy: 0.3593 - val_loss: 1.8150 - val_accuracy: 0.6330
Epoch 3/20
63/63 [=====] - 1s 15ms/step - loss: 1.4413 - accuracy: 0.6844 - val_loss: 1.4487 - val_accuracy: 0.6650
Epoch 4/20
63/63 [=====] - 1s 20ms/step - loss: 1.1981 - accuracy: 0.7129 - val_loss: 1.3661 - val_accuracy: 0.6760
Epoch 5/20
63/63 [=====] - 2s 24ms/step - loss: 1.0759 - accuracy: 0.7378 - val_loss: 1.3343 - val_accuracy: 0.6880
Epoch 6/20
63/63 [=====] - 1s 22ms/step - loss: 0.9850 - accuracy: 0.7600 - val_loss: 1.3316 - val_accuracy: 0.7070
Epoch 7/20
63/63 [=====] - 1s 15ms/step - loss: 0.9110 - accuracy: 0.7767 - val_loss: 1.3018 - val_accuracy: 0.7060
Epoch 8/20
63/63 [=====] - 1s 16ms/step - loss: 0.8501 - accuracy: 0.7955 - val_loss: 1.3672 - val_accuracy: 0.7080
Epoch 9/20
63/63 [=====] - 1s 19ms/step - loss: 0.7946 - accuracy: 0.8067 - val_loss: 1.3694 - val_accuracy: 0.7070
Epoch 10/20
63/63 [=====] - 1s 19ms/step - loss: 0.7482 - accuracy: 0.8175 - val_loss: 1.3640 - val_accuracy: 0.7140
```

```
[23] y_test = to_categorical(test_labels)
      results = model.evaluate(x_test, y_test)
      print(results)

71/71 [=====] - 0s 4ms/step - loss: 1.9224 - accuracy: 0.6919
[1.9223910570144653, 0.6918966770172119]
```

The model performs considerably worse with a much lower accuracy than before due to information bottleneck. This happens because we are trying to compress a lot of information into an intermediate space that is too low dimensional (4 units) and then attempting to recover information from this that describes 46 classes i.e. we're trying to recover the separation hyperplanes for 46 classes. The model performance is around 70% accuracy which means the model is actually able to cram most of the important information into these 4-dimensional representations but clearly not all of it since it resulted in a lower accuracy than before. Consider this when you are designing your own neural network architecture.

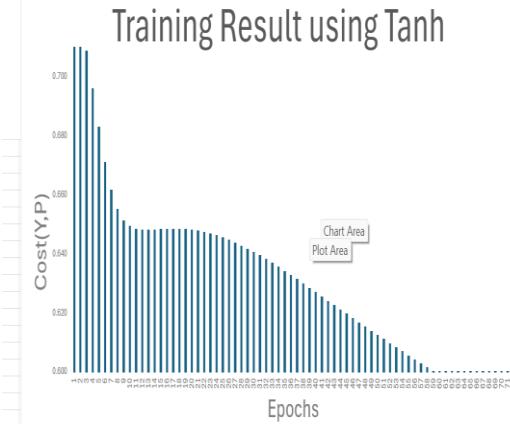
## VII. WEEK 5

**A. Resource:** With assumption that you have access to Microsoft Excel, create a new template with the following setup. From there, you start implement the multilayer perceptron logistic regression for solving XOR function approximation based on the formulas discussed in the lecture.

		Using Sigmoid Activation Function																																
Epoch	X	Y	V1	V2	b1	b2	L <sub>ReLU</sub>	Z1	A1 + sigmoid(Z1)	22	P <sup>T</sup> A2 + b2Z2	C	d <sub>0</sub> b2	d <sub>0</sub> d <sub>2</sub>	d <sub>0</sub> d <sub>3</sub>	d <sub>0</sub> d <sub>4</sub>	d <sub>0</sub> d <sub>5</sub>	V1	b1	v2	b2													
1	0	0	-0.7	-0.7	-0.5	0.0	0.0	0.0	-0.7	-0.7	-14	0.60	0.332	0.322	0.186	0.4	0.599	0.722	-0.098	0.009	0.001	-0.003	-0.004	-0.03	-0.631	-0.689	0.02	-0.643	-0.403					
1	0	1	-0.5	0.7	0.4	0.0	0.0	0.0	-0.5	-0.5	0.2	0.500	0.650	0.738	0.560	0.767	0.689	-0.083	0.007	-0.009	0.07	0.007	0.005	0.09	-0.522	0.872	-0.044	0.178						
1	1	0	0.1	0.7	0.9	0.0	0.0	0.0	-0.1	0.6	0.6	0.500	0.688	0.475	0.546	0.417	0.620	0.000	-0.001	-0.002	0.034	0.012	0.008	0.028	-0.138	0.841	-0.068	0.680						
1	1	1	0.5	0.5	0.5	0.0	0.0	0.0	0.5	0.5	0.5	0.500	0.650	0.650	0.560	0.417	0.620	0.000	-0.001	-0.002	0.034	0.012	0.008	0.028	-0.138	0.841	-0.068	0.680						
2	0	0	0	-0.893	-0.689	-0.642	0.058	-0.403	3	0.0	-0.7	-0.7	-13	0.569	0.342	0.341	0.207	-0.3	0.419	0.897	-0.058	-0.014	-0.045	-0.07	0.019	0.020	-0.013	0.009	0.006	-0.78	-0.768	0.008	-0.599	-0.263
2	0	1	1	-0.522	0.672	0.78	-0.044	0.0	0.0	0.6	-0.6	0.1	0.499	0.652	0.362	0.527	0.417	0.620	0.000	-0.001	-0.021	0.034	0.012	0.008	0.028	-0.138	0.841	-0.068	0.680					
2	1	0	0	-0.128	0.646	0.869	-0.088	0.0	-0.1	0.6	-0.2	0.4	0.478	0.936	0.444	0.604	0.417	0.620	0.000	-0.001	-0.022	0.034	0.012	0.008	0.006	-0.008	-0.125	0.841	-0.068	0.729				
2	1	1	0	0	-0.522	0.672	0.78	-0.044	0.0	0.0	0.6	0.1	0.499	0.652	0.362	0.527	0.417	0.620	0.000	-0.001	-0.021	0.034	0.012	0.008	0.028	-0.138	0.841	-0.068	0.680					
3	0	0	0	-0.78	-0.78	-0.618	0.004	-0.304	3	0.0	-0.7	-0.7	-14	0.569	0.329	0.338	0.192	-0.2	0.457	0.892	0.114	0.000	-0.008	-0.08	0.008	-0.013	0.005	0.004	0.004	-0.744	-0.724	-0.016	-0.618	-0.304
3	0	1	1	-0.520	0.672	0.241	-0.044	0.0	0.0	0.6	-0.6	0.1	0.490	0.653	0.363	0.527	0.417	0.620	0.000	-0.001	-0.007	0.007	0.008	0.005	0.002	0.002	-0.524	0.884	0.047	0.216				
3	1	0	0	-0.122	0.646	0.725	-0.068	0.0	-0.1	0.6	-0.2	0.4	0.443	0.841	0.452	0.512	0.417	0.620	0.000	-0.001	-0.019	0.024	0.008	0.005	0.003	-0.125	0.841	-0.068	0.680					
3	1	1	0	0	-0.520	0.672	0.241	-0.044	0.0	0.0	0.6	0.1	0.490	0.653	0.363	0.527	0.417	0.620	0.000	-0.001	-0.007	0.007	0.008	0.005	0.002	0.002	-0.524	0.884	0.047	0.216				
4	0	0	0	-0.78	-0.78	-0.618	0.004	-0.304	3	0.0	-0.7	-0.7	-14	0.569	0.329	0.338	0.192	-0.2	0.457	0.892	0.114	0.000	-0.008	-0.08	0.008	-0.013	0.005	0.004	-0.744	-0.724	-0.016	-0.618	-0.304	
4	0	1	1	-0.524	0.683	0.276	-0.047	0.0	0.0	0.6	-0.6	0.1	0.486	0.656	0.366	0.524	0.417	0.620	0.000	-0.001	-0.008	0.008	0.005	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
4	1	0	0	-0.122	0.646	0.725	-0.068	0.0	-0.1	0.6	-0.2	0.4	0.443	0.841	0.452	0.512	0.417	0.620	0.000	-0.001	-0.019	0.024	0.008	0.005	0.003	-0.125	0.841	-0.068	0.680					
4	1	1	0	0	-0.524	0.683	0.276	-0.047	0.0	0.0	0.6	0.1	0.486	0.656	0.366	0.524	0.417	0.620	0.000	-0.001	-0.008	0.008	0.005	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
5	0	0	0	-0.744	-0.724	-0.618	-0.016	-0.287	3	0.0	-0.7	-0.7	-15	0.497	0.324	0.330	0.198	-0.2	0.463	0.892	0.118	0.002	-0.008	-0.08	0.007	-0.012	0.005	0.004	-0.753	-0.725	0.023	-0.632	-0.286	
5	0	1	1	-0.520	0.683	0.276	-0.047	0.0	0.0	0.6	-0.6	0.1	0.486	0.656	0.366	0.524	0.417	0.620	0.000	-0.001	-0.008	0.008	0.005	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
5	1	0	0	-0.122	0.646	0.725	-0.068	0.0	-0.1	0.6	-0.2	0.4	0.443	0.841	0.452	0.512	0.417	0.620	0.000	-0.001	-0.019	0.024	0.008	0.005	0.003	-0.125	0.841	-0.068	0.680					
5	1	1	0	0	-0.520	0.683	0.276	-0.047	0.0	0.0	0.6	0.1	0.486	0.656	0.366	0.524	0.417	0.620	0.000	-0.001	-0.008	0.008	0.005	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
6	0	0	0	-0.709	-0.728	-0.625	-0.020	-0.286	3	0.0	-0.8	-0.8	-15	0.497	0.285	0.349	0.180	-0.2	0.463	0.892	0.120	0.002	-0.009	-0.09	0.005	0.004	0.004	-0.778	-0.747	-0.037	-0.627	-0.286		
6	0	1	1	-0.522	0.685	0.278	-0.043	0.0	0.0	0.6	-0.6	0.1	0.487	0.657	0.368	0.526	0.417	0.620	0.000	-0.001	-0.006	0.006	0.007	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
6	1	0	0	-0.122	0.646	0.868	-0.078	0.0	-0.1	0.6	-0.2	0.4	0.443	0.841	0.450	0.510	0.417	0.620	0.000	-0.001	-0.019	0.024	0.008	0.005	0.003	-0.125	0.841	-0.068	0.680					
6	1	1	0	0	-0.522	0.685	0.278	-0.043	0.0	0.0	0.6	0.1	0.487	0.657	0.368	0.526	0.417	0.620	0.000	-0.001	-0.006	0.006	0.007	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
7	0	0	0	-0.775	-0.743	-0.627	-0.033	-0.286	3	0.0	-0.9	-0.9	-16	0.497	0.285	0.349	0.175	-0.2	0.463	0.892	0.120	0.002	-0.009	-0.09	0.005	0.004	0.004	-0.791	-0.754	-0.043	-0.632	-0.276		
7	0	1	1	-0.523	0.695	0.298	-0.056	0.0	0.0	0.6	-0.6	0.1	0.486	0.658	0.368	0.525	0.417	0.620	0.000	-0.001	-0.008	0.008	0.007	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
7	1	0	0	-0.122	0.646	0.868	-0.078	0.0	-0.1	0.6	-0.2	0.4	0.440	0.840	0.450	0.510	0.417	0.620	0.000	-0.001	-0.019	0.024	0.008	0.005	0.003	-0.125	0.841	-0.068	0.680					
7	1	1	0	0	-0.523	0.695	0.298	-0.056	0.0	0.0	0.6	0.1	0.486	0.658	0.368	0.525	0.417	0.620	0.000	-0.001	-0.008	0.008	0.007	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
8	0	0	0	-0.781	-0.750	-0.623	-0.042	-0.275	3	0.0	-0.8	-0.8	-18	0.490	0.280	0.343	0.163	-0.2	0.461	0.891	0.118	0.002	-0.009	-0.09	0.005	0.004	0.004	-0.808	-0.777	-0.062	-0.641	-0.272		
8	0	1	1	-0.519	0.695	0.298	-0.056	0.0	0.0	0.6	-0.6	0.1	0.485	0.652	0.362	0.516	0.417	0.620	0.000	-0.001	-0.008	0.008	0.007	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
8	1	0	0	-0.121	0.646	0.868	-0.080	0.0	-0.1	0.6	-0.2	0.4	0.440	0.840	0.450	0.510	0.417	0.620	0.000	-0.001	-0.019	0.024	0.008	0.005	0.003	-0.125	0.841	-0.068	0.680					
8	1	1	0	0	-0.519	0.695	0.298	-0.056	0.0	0.0	0.6	0.1	0.485	0.652	0.362	0.516	0.417	0.620	0.000	-0.001	-0.008	0.008	0.007	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
9	0	0	0	-0.898	-0.772	-0.644	-0.062	-0.272	3	-0.1	-0.8	-0.9	-16	0.497	0.269	0.327	0.163	-0.2	0.460	0.891	0.118	0.003	-0.009	-0.08	0.007	-0.012	0.005	0.004	0.004	-0.828	-0.797	-0.049	-0.649	-0.267
9	0	1	1	-0.523	0.647	0.202	-0.062	0.0	0.0	0.6	-0.6	0.1	0.485	0.655	0.355	0.507	0.417	0.620	0.000	-0.001	-0.008	0.008	0.007	0.006	0.003	0.003	-0.525	0.885	0.033	0.219				
9	1	0	0	-0.121	0.646	0.868	-0.070	0.0	-0.1	0.6	-0.2	0.4	0.440	0.840	0.450	0.507	0.417	0.620	0.000	-0.001	-0.019	0.024	0.008	0.005	0.003	-0.125	0.841	-0.068	0.680					
9	1																																	

2. Replace Sigmoid function with Tanh function. Explain your finding along with your analysis evidence.

plot the Cost function value at every epoch, up to 70 epochs.



## **Answer-**

## Impact of Activation Functions on Training:

This experiment compared the performance of two common activation functions, tanh and sigmoid, in the context of training our time series forecasting model. The initial cost (loss) value was lower for the sigmoid function (0.732) compared to the tanh function (0.771). While both functions converged to similar cost values by the end of training, the tanh function exhibited slightly higher cost throughout the middle epochs. Additionally, the cost reduction with tanh displayed more variability compared to the smoother decline observed with sigmoid.

These observations suggest potential advantages and disadvantages for each activation function:

- **Tanh:** May offer a faster initial reduction in cost, potentially leading to quicker convergence in the early stages of training.
  - **Sigmoid:** Provides a smoother and potentially more stable cost reduction throughout training.

Therefore, the selection of the activation function might depend on the specific training goals. If faster initial convergence is prioritized, tanh could be a viable choice. However, if overall training stability is more crucial, sigmoid might be the preferred option.

## Week 5 Workshop

- B. Please complete the following notebook and the attached workshop, record the outcome in your log book and demonstrate your understand to your best level.

<https://colab.research.google.com/drive/1ggISnRkkSNBnYb8sxqyIRBgmq3PulpmG?usp=sharing>

### House Price Prediction with Deep Learning

- **Dataset:** The Boston House Prices dataset was employed for this project. This dataset consists of 506 data points with 13 features each, representing various characteristics of houses. The data was normalized to ensure all features are on a similar scale for effective model training.

The screenshot shows a Jupyter Notebook interface. On the left, there's a sidebar with sections like 'Regression Example in Tensorflow : Predicting house prices' and 'The Boston Housing Price dataset'. The main area contains Python code for loading the dataset and normalizing the data. On the right, there's a section titled 'Preparing the data' with sub-sections for 'Normalizing the data' and 'Note: the quantities for normalising the test data comes from the training data as shown below.'

```
[1]: import numpy as np
import tensorflow as tf
import keras
from tensorflow import keras
from tensorflow.keras import layers
from tensorflow.keras.datasets import boston_housing
tf.random.set_seed(123)
np.random.seed(123)
(train_data, train_targets), (test_data, test_targets) = boston_housing.load_data()

[6]: mean = train_data.mean(axis=0)
train_data -= mean
std = train_data.std(axis=0)
train_data /= std
test_data -= mean
test_data /= std
```

- **Model Architecture:** A deep learning model with five dense layers was constructed. The first layer has **13 neurons** corresponding to the number of input features. The subsequent layers have **64 neurons** each, designed to capture complex relationships between features and house prices.

The screenshot shows a Jupyter Notebook interface. It displays a Python function 'build\_model()' which creates a sequential model with five layers. Below the function, there's a section titled 'K-fold validation' with explanatory text about the process. A scroll bar is visible on the right side of the code cell.

```
def build_model():
    model = keras.Sequential([
        layers.Dense(13, activation="elu"),
        layers.Dense(64, activation="elu"),
        layers.Dense(64, activation="elu"),
        layers.Dense(64, activation="elu"),
        layers.Dense(1)
    ])
    model.compile(optimizer="rmsprop", loss="mse", metrics=["mae"])
    return model
```

Dividing the Dataset: First, the dataset is divided into k equal-sized subsets or folds.  
Training and Testing: The cross-validation process then iterates through k rounds. In each round, one of the k subsets is used as the test set, while the remaining k-1 subsets are used as the training set.  
Model Training and Evaluation: The model is trained on the training set and evaluated on the test set. This process is repeated k times, each time using a different subset as the test set.  
Performance Evaluation: Once all rounds are completed, the performance metrics (such as accuracy, precision, recall, etc.) from each iteration are averaged to provide a overall measure of the model's performance.  
Final Model: After cross-validation, if the performance is satisfactory, the model can be trained on the entire dataset using the selected hyperparameters.

- **K-Fold Cross-Validation:** To evaluate model performance and mitigate overfitting, 4-fold cross-validation was implemented. This technique divides the data into four folds, trains the model on three folds while holding out the remaining fold for validation. This process is repeated four times, ensuring all data points are used for both training and validation.

```
[24] k = 4
num_val_samples = len(train_data) // k
num_epochs = 100
all_scores = []
for i in range(k):
    print(f"Processing fold #{i}")
    val_data = train_data[i * num_val_samples: (i + 1) * num_val_samples]
    val_targets = train_targets[i * num_val_samples: (i + 1) * num_val_samples]
    partial_train_data = np.concatenate(
        [train_data[:i * num_val_samples],
         train_data[(i + 1) * num_val_samples:]]),
    axis=0)
    partial_train_targets = np.concatenate(
        [train_targets[:i * num_val_samples],
         train_targets[(i + 1) * num_val_samples:]]),
    axis=0)
    model = build_model()
    model.fit(partial_train_data, partial_train_targets,
              epochs=num_epochs, batch_size=16, verbose=0)
    val_mse, val_mae = model.evaluate(val_data, val_targets, verbose=0)
    all_scores.append(val_mae)

⇒ Processing fold #0
Processing fold #1
Processing fold #2
Processing fold #3

Observe the scores from each of the k-fold runs. We're using MAE, so the lower is better. For example, for the first score, this means that across all predictions, your model is off by $1966.4
```

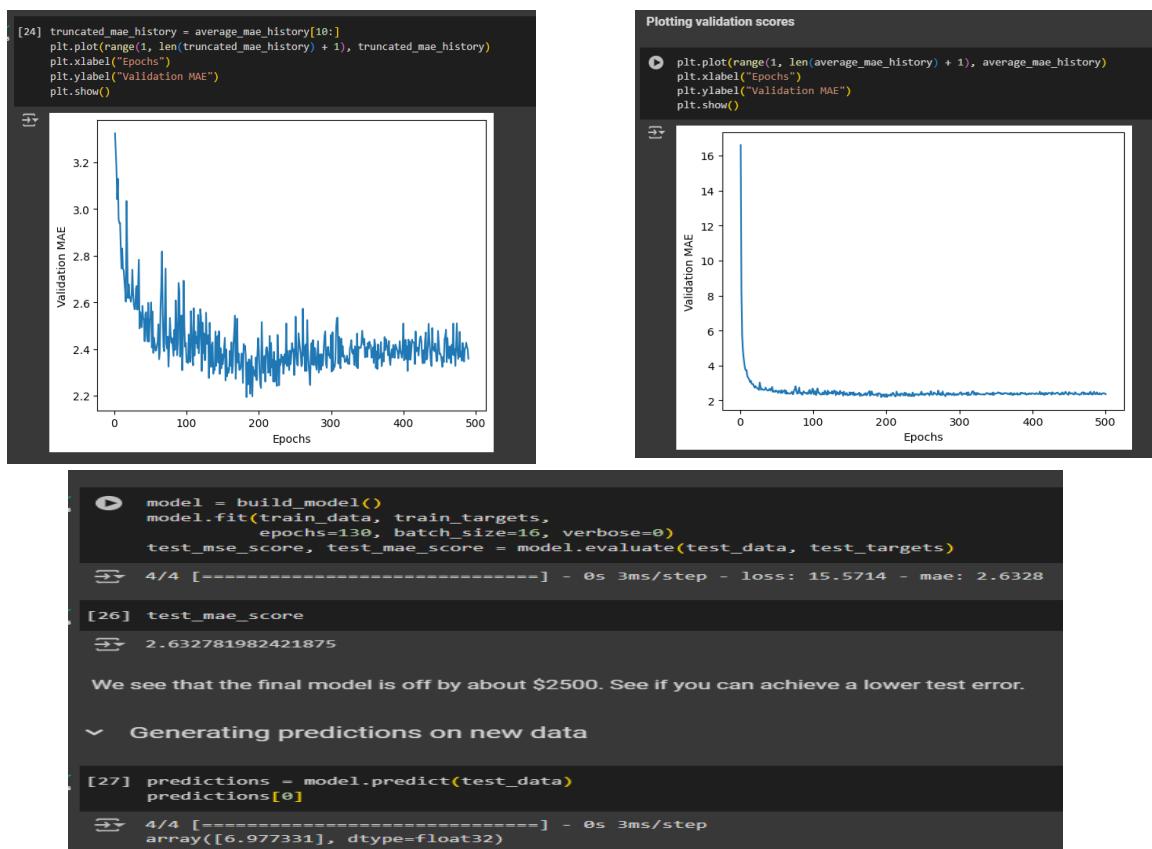
[26] all\_scores

```
⇒ [2.292691469192505, 2.453836679458618, 2.5215258598327637, 2.378781795501709]
```

We then take the average of the MAE scores for all folds:

```
⇒ np.mean(all_scores)
⇒ 2.411700950996399
```

- **Initial Training (100 Epochs and 500 epochs):** The model was initially trained for 100 epochs using K-fold cross-validation. The results from each fold were combined to obtain a mean Mean Absolute Error (MAE) of **2.41** on the validation sets. This score provides a baseline performance metric.
- **Hyperparameter Tuning and Training (130 Epochs):** To potentially improve performance, we further investigated the model's hyperparameters. Based on this analysis, the model was then trained for a longer duration (130 epochs) with the final chosen hyperparameters. However, this final training resulted in a higher MAE of **2.63**. We can visualize the trend of the MAE during training in the included plot



## VIII. WEEK 6

A. Please complete the following notebook, record the outcome in your log book and demonstrate your understand to your best level.

<https://colab.research.google.com/drive/1tgvOaNuQ6aC4zIQpc1H8xLdr0cbP6Ly7?usp=sharing>

The screenshot shows two code cells from a Jupyter Notebook. The first cell adds white-noise channels or all-zero channels to the MNIST dataset:

```
from tensorflow.keras.datasets import mnist
import numpy as np

(train_images, train_labels), _ = mnist.load_data()
train_images = train_images.reshape((60000, 28 * 28))
train_images = train_images.astype("float32") / 255

train_images_with_noise_channels = np.concatenate(
    [train_images, np.random.randint(0, 255, (len(train_images), 784))], axis=1)

train_images_with_zeros_channels = np.concatenate(
    [train_images, np.zeros((len(train_images), 784))], axis=1)
```

The second cell shows the lengths of the modified datasets:

```
[2] len(train_images_with_noise_channels)
[3] len(train_labels)
```

The third cell shows the shape of the noisy dataset:

```
[4] train_images_with_noise_channels.shape
```

The second part of the notebook contains code for training a neural network on the noisy and zero-channel datasets:

```
from tensorflow import keras
from tensorflow.keras import layers
import tensorflow as tf

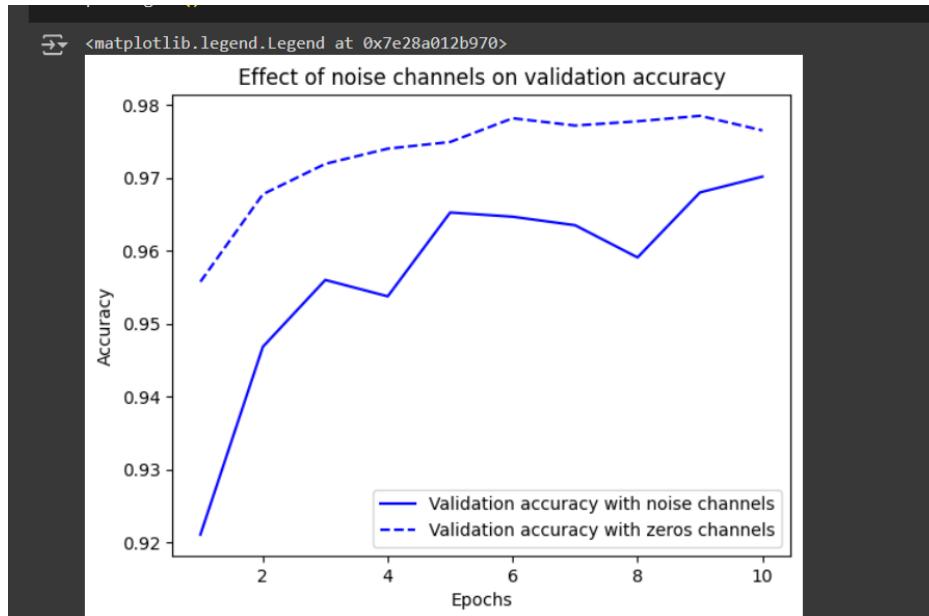
def get_model():
    model = keras.Sequential([
        layers.Dense(512, activation="relu"),
        layers.Dense(10, activation="softmax")
    ])
    model.compile(optimizer="rmsprop",
                  loss="sparse_categorical_crossentropy",
                  metrics=["accuracy"])
    return model

model = get_model()
history_noise = model.fit(
    train_images_with_noise_channels, train_labels,
    epochs=10,
    batch_size=128,
    validation_split=0.2)

model = get_model()
history_zeros = model.fit(
    train_images_with_zeros_channels, train_labels,
    epochs=10,
    batch_size=128,
    validation_split=0.2)
```

### Plotting a validation accuracy comparison

```
0s   import matplotlib.pyplot as plt
     val_acc_noise = history_noise.history["val_accuracy"]
     val_acc_zeros = history_zeros.history["val_accuracy"]
     epochs = range(1, 11)
     plt.plot(epochs, val_acc_noise, "b-",
              label="Validation accuracy with noise channels")
     plt.plot(epochs, val_acc_zeros, "b--",
              label="Validation accuracy with zeros channels")
     plt.title("Effect of noise channels on validation accuracy")
     plt.xlabel("Epochs")
     plt.ylabel("Accuracy")
     plt.legend()
```



**Answer -**Our investigation into the impact of data noise on model performance revealed a clear trend. Two datasets were employed: one containing noise and one without. Both datasets were trained on a model featuring two dense layers with ReLU and sigmoid activation functions. The RMSProp optimizer was utilized, and training proceeded for 10 epochs. Subsequent analysis of validation accuracy showed that the dataset absent of noise demonstrably achieved superior validation accuracy compared to the dataset containing noise. This finding underscores the importance of data quality in machine learning, as noise can significantly impede model performance.

## Now we have build couple of models for experiments -

### Training a MNIST model with an incorrectly high learning rate

Let's set the learning rate = 1 in the below cell. We can see that this model reaches a low training and validation accuracy and cannot get past that.

```
▶ (train_images, train_labels), _ = mnist.load_data()
train_images = train_images.reshape((60000, 28 * 28))
train_images = train_images.astype("float32") / 255

model = keras.Sequential([
    layers.Dense(512, activation="relu"),
    layers.Dense(10, activation="softmax")
])
model.compile(optimizer=keras.optimizers.RMSprop(1.),
              loss="sparse_categorical_crossentropy",
              metrics=["accuracy"])
model.fit(train_images, train_labels,
          epochs=10,
          batch_size=128,
          validation_split=0.2)

→ Epoch 1/10
375/375 [=====] - 2s 4ms/step - loss: 695.2733 - accuracy: 0.3734 - val_loss: 2.2379 - val_accuracy: 0.2432
Epoch 2/10
375/375 [=====] - 1s 3ms/step - loss: 2.9671 - accuracy: 0.2527 - val_loss: 2.9329 - val_accuracy: 0.2317
Epoch 3/10
375/375 [=====] - 1s 3ms/step - loss: 2.4130 - accuracy: 0.2141 - val_loss: 2.3540 - val_accuracy: 0.2476
Epoch 4/10
375/375 [=====] - 1s 3ms/step - loss: 2.6342 - accuracy: 0.2092 - val_loss: 2.1454 - val_accuracy: 0.1928
```

## 1. Model with High learning rate – This was absolutely disaster we cannot reach to global minima we reached only 0.21 validation accuracy.

### The same model with a more appropriate learning rate

```
▶ model = keras.Sequential([
    layers.Dense(512, activation="relu"),
    layers.Dense(10, activation="softmax")
])
model.compile(optimizer=keras.optimizers.RMSprop(1e-2),
              loss="sparse_categorical_crossentropy",
              metrics=["accuracy"])
model.fit(train_images, train_labels,
          epochs=10,
          batch_size=128,
          validation_split=0.2)

→ Epoch 1/10
375/375 [=====] - 2s 4ms/step - loss: 0.3629 - accuracy: 0.9118 - val_loss: 0.1518 - val_accuracy: 0.9548
Epoch 2/10
375/375 [=====] - 2s 4ms/step - loss: 0.1284 - accuracy: 0.9638 - val_loss: 0.1752 - val_accuracy: 0.9564
Epoch 3/10
375/375 [=====] - 1s 4ms/step - loss: 0.0975 - accuracy: 0.9737 - val_loss: 0.1548 - val_accuracy: 0.9679
Epoch 4/10
375/375 [=====] - 1s 3ms/step - loss: 0.0818 - accuracy: 0.9792 - val_loss: 0.1549 - val_accuracy: 0.9701
Epoch 5/10
375/375 [=====] - 1s 3ms/step - loss: 0.0688 - accuracy: 0.9824 - val_loss: 0.2219 - val_accuracy: 0.9625
Epoch 6/10
375/375 [=====] - 1s 3ms/step - loss: 0.0634 - accuracy: 0.9849 - val_loss: 0.2009 - val_accuracy: 0.9661
```

## 2. Model with appropriate learning rate- This model perform very well we received 0.976 validation accuracy

#### A simple logistic regression on MNIST

This model contains one layer.

```
[9] model = keras.Sequential([layers.Dense(10, activation="softmax")])
    model.compile(optimizer="rmsprop",
                  loss="sparse_categorical_crossentropy",
                  metrics=["accuracy"])
history_small_model = model.fit(
    train_images, train_labels,
    epochs=20,
    batch_size=128,
    validation_split=0.2)

→ Epoch 1/20
375/375 [=====] - 2s 4ms/step - loss: 0.6652 - accuracy: 0.8373 - val_loss: 0.3627 - val_accuracy: 0.9016
Epoch 2/20
375/375 [=====] - 1s 4ms/step - loss: 0.3532 - accuracy: 0.9029 - val_loss: 0.3100 - val_accuracy: 0.9147
Epoch 3/20
375/375 [=====] - 1s 3ms/step - loss: 0.3178 - accuracy: 0.9110 - val_loss: 0.2938 - val_accuracy: 0.9197
Epoch 4/20
375/375 [=====] - 1s 3ms/step - loss: 0.3018 - accuracy: 0.9155 - val_loss: 0.2839 - val_accuracy: 0.9219
Epoch 5/20
375/375 [=====] - 1s 3ms/step - loss: 0.2924 - accuracy: 0.9185 - val_loss: 0.2793 - val_accuracy: 0.9219
Epoch 6/20
```

### 3.Simple model with one layer – Validation accuracy we get around 0.92 which is decent model.

Lets increase the capacity of the model and evaluate its performance compared to the model with a single layer.

```
[11] model = keras.Sequential([
    layers.Dense(96, activation="relu"),
    layers.Dense(96, activation="relu"),
    layers.Dense(10, activation="softmax"),
])
model.compile(optimizer="rmsprop",
              loss="sparse_categorical_crossentropy",
              metrics=["accuracy"])
history_large_model = model.fit(
    train_images, train_labels,
    epochs=20,
    batch_size=128,
    validation_split=0.2)

→ Epoch 1/20
375/375 [=====] - 2s 4ms/step - loss: 0.3657 - accuracy: 0.8956 - val_loss: 0.1905 - val_accuracy: 0.9454
Epoch 2/20
375/375 [=====] - 1s 3ms/step - loss: 0.1593 - accuracy: 0.9528 - val_loss: 0.1427 - val_accuracy: 0.9588
Epoch 3/20
375/375 [=====] - 2s 4ms/step - loss: 0.1111 - accuracy: 0.9668 - val_loss: 0.1103 - val_accuracy: 0.9668
Epoch 4/20
375/375 [=====] - 2s 4ms/step - loss: 0.0852 - accuracy: 0.9739 - val_loss: 0.0997 - val_accuracy: 0.9704
Epoch 5/20
375/375 [=====] - 1s 3ms/step - loss: 0.0687 - accuracy: 0.9791 - val_loss: 0.0987 - val_accuracy: 0.9716
Epoch 6/20
```

### 4. 3 layer model - this model perform very well. Get around 0.97 validation accuracy but tends to overfit after 8 epochs

## Now we have comparing three model – Original model , Less neuron model,High neuron layer

### Original model

```

from tensorflow.keras.datasets import imdb
(train_data, train_labels), _ = imdb.load_data(num_words=10000)

def vectorize_sequences(sequences, dimension=10000):
    results = np.zeros((len(sequences), dimension))
    for i, sequence in enumerate(sequences):
        results[i, sequence] = 1.
    return results
train_data = vectorize_sequences(train_data)

model = keras.Sequential([
    layers.Dense(16, activation="relu"),
    layers.Dense(16, activation="relu"),
    layers.Dense(1, activation="sigmoid")
])
model.compile(optimizer="rmsprop",
              loss="binary_crossentropy",
              metrics=["accuracy"])
history_original = model.fit(train_data, train_labels,
                             epochs=20, batch_size=512, validation_split=0.4)

```

Downloading data from <https://storage.googleapis.com/tensorflow/tf-keras-datasets/imdb.npz>  
17464789/17464789 [=====] - 0s 0us/step  
Epoch 1/20  
30/30 [=====] - 4s 65ms/step - loss: 0.5412 - accuracy: 0.7629 - val\_loss: 0.3932 - val\_accuracy: 0.8695  
Epoch 2/20  
30/30 [=====] - 1s 25ms/step - loss: 0.3254 - accuracy: 0.8903 - val\_loss: 0.3089 - val\_accuracy: 0.8842  
Epoch 3/20  
30/30 [=====] - 1s 30ms/step - loss: 0.2417 - accuracy: 0.9193 - val\_loss: 0.2826 - val\_accuracy: 0.8900  
Epoch 4/20  
30/30 [=====] - 1s 30ms/step - loss: 0.1928 - accuracy: 0.9373 - val\_loss: 0.2913 - val\_accuracy: 0.8831  
Epoch 5/20

### Version of the model with lower capacity

```

model = keras.Sequential([
    layers.Dense(4, activation="relu"),
    layers.Dense(4, activation="relu"),
    layers.Dense(1, activation="sigmoid")
])
model.compile(optimizer="rmsprop",
              loss="binary_crossentropy",
              metrics=["accuracy"])
history_small_model = model.fit(
    train_data, train_labels,
    epochs=20, batch_size=512, validation_split=0.4)

```

Epoch 1/20  
30/30 [=====] - 3s 62ms/step - loss: 0.6320 - accuracy: 0.7449 - val\_loss: 0.5654 - val\_accuracy: 0.8539  
Epoch 2/20  
30/30 [=====] - 1s 22ms/step - loss: 0.5043 - accuracy: 0.8711 - val\_loss: 0.4626 - val\_accuracy: 0.8692  
Epoch 3/20  
30/30 [=====] - 1s 24ms/step - loss: 0.4065 - accuracy: 0.8911 - val\_loss: 0.3919 - val\_accuracy: 0.8732  
Epoch 4/20  
30/30 [=====] - 1s 22ms/step - loss: 0.3348 - accuracy: 0.9038 - val\_loss: 0.3482 - val\_accuracy: 0.8763  
Epoch 5/20  
30/30 [=====] - 1s 25ms/step - loss: 0.2820 - accuracy: 0.9172 - val\_loss: 0.3204 - val\_accuracy: 0.8784  
Epoch 6/20  
30/30 [=====] - 1s 21ms/step - loss: 0.2425 - accuracy: 0.9259 - val\_loss: 0.2910 - val\_accuracy: 0.8883  
Epoch 7/20  
30/30 [=====] - 1s 24ms/step - loss: 0.2121 - accuracy: 0.9325 - val\_loss: 0.2815 - val\_accuracy: 0.8906  
Epoch 8/20  
30/30 [=====] - 1s 22ms/step - loss: 0.1885 - accuracy: 0.9397 - val\_loss: 0.2750 - val\_accuracy: 0.8934  
Epoch 9/20

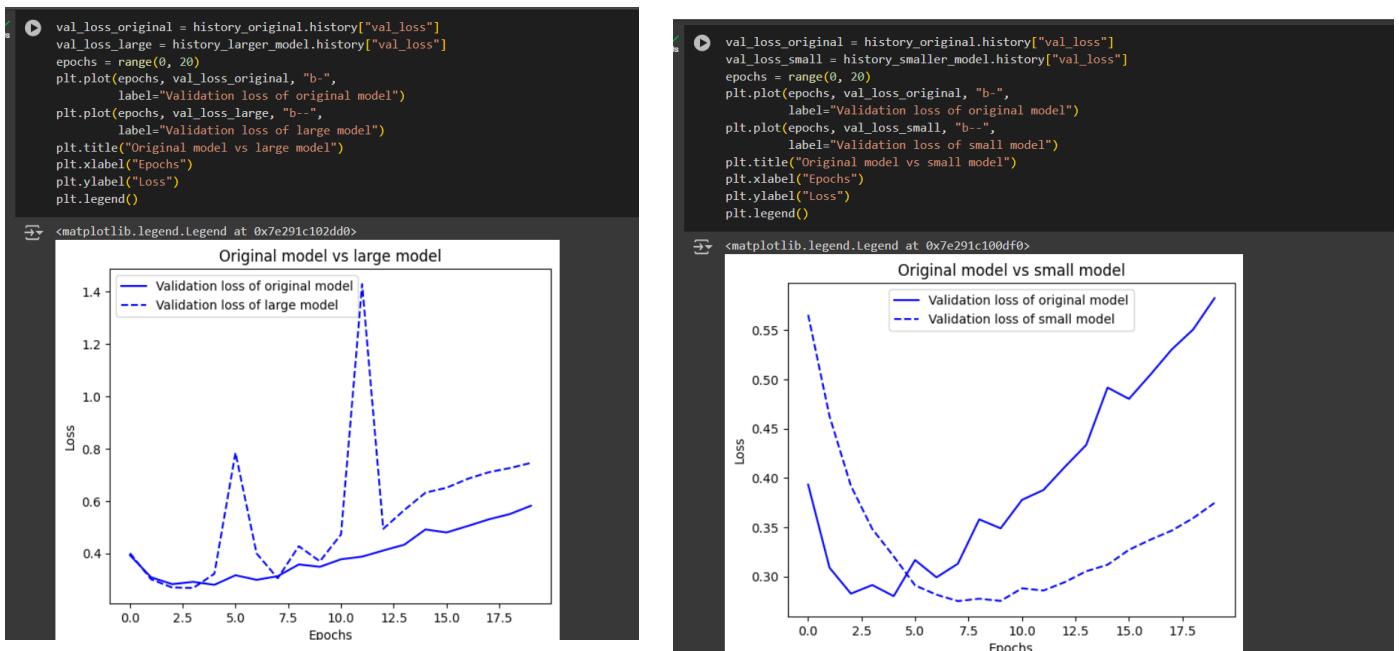
### Version of the model with higher capacity

```

[15] model = keras.Sequential([
    layers.Dense(512, activation="relu"),
    layers.Dense(512, activation="relu"),
    layers.Dense(1, activation="sigmoid")
])
model.compile(optimizer="rmsprop",
              loss="binary_crossentropy",
              metrics=["accuracy"])
history_larger_model = model.fit(
    train_data, train_labels,
    epochs=20, batch_size=512, validation_split=0.4)

```

Epoch 1/20  
30/30 [=====] - 3s 74ms/step - loss: 0.5674 - accuracy: 0.7121 - val\_loss: 0.4010 - val\_accuracy: 0.8248  
Epoch 2/20  
30/30 [=====] - 1s 26ms/step - loss: 0.3255 - accuracy: 0.8662 - val\_loss: 0.3013 - val\_accuracy: 0.8762  
Epoch 3/20  
30/30 [=====] - 1s 25ms/step - loss: 0.2456 - accuracy: 0.8987 - val\_loss: 0.2698 - val\_accuracy: 0.8894  
Epoch 4/20  
30/30 [=====] - 1s 25ms/step - loss: 0.1799 - accuracy: 0.9323 - val\_loss: 0.2682 - val\_accuracy: 0.8923  
Epoch 5/20  
30/30 [=====] - 1s 25ms/step - loss: 0.1314 - accuracy: 0.9527 - val\_loss: 0.3224 - val\_accuracy: 0.8799  
Epoch 6/20  
30/30 [=====] - 1s 26ms/step - loss: 0.0998 - accuracy: 0.9668 - val\_loss: 0.7850 - val\_accuracy: 0.7730  
Epoch 7/20



The impact of the number of neurons on model training was explored by training models with varying neuron counts in each layer. We analyzed the resulting graphs (refer to Figure above) and identified several key observations.

- **Original vs. Large Model:** Compared to the original model, the large model with a higher number of neurons exhibited signs of inconsistency and a tendency to overfit readily.
- **Original vs. Small Model:** The small model, with a lower number of neurons per layer, demonstrated slower convergence. While it eventually reached a minimum point around 8-10 epochs, it also displayed overfitting tendencies thereafter.

Now we add some penalties to over existing models

1. Apply L2 regularization
2. Apply Dropouts

With this it help to stop overfitting problem in our model.

Let compare both technique with our model

```
✓ Adding weight regularization

Adding L2 weight regularization to the model

from tensorflow.keras import regularizers
model = keras.Sequential([
    layers.Dense(16,
        kernel_regularizer=regularizers.l2(0.002),
        activation="relu"),
    layers.Dense(16,
        kernel_regularizer=regularizers.l2(0.002),
        activation="relu"),
    layers.Dense(1, activation="sigmoid")
])
model.compile(optimizer="rmsprop",
    loss="binary_crossentropy",
    metrics=["accuracy"])
history_l2_reg = model.fit(
    train_data, train_labels,
    epochs=20, batch_size=512, validation_split=0.4)

Epoch 1/20
30/30 [=====] - 3s 85ms/step - loss: 0.6353 - accuracy: 0.7647 - val_loss: 0.5002 - val_accuracy: 0.8656
Epoch 2/20
30/30 [=====] - 1s 25ms/step - loss: 0.4284 - accuracy: 0.8908 - val_loss: 0.4127 - val_accuracy: 0.8729
Epoch 3/20
30/30 [=====] - 1s 25ms/step - loss: 0.3451 - accuracy: 0.9099 - val_loss: 0.3654 - val_accuracy: 0.8881
Epoch 4/20
30/30 [=====] - 1s 24ms/step - loss: 0.3003 - accuracy: 0.9214 - val_loss: 0.3516 - val_accuracy: 0.8897
Epoch 5/20
30/30 [=====] - 1s 21ms/step - loss: 0.2775 - accuracy: 0.9293 - val_loss: 0.3783 - val_accuracy: 0.8720
Epoch 6/20
```

```
Adding more layers and changing the number of units in each layer:

from tensorflow.keras import regularizers
from tensorflow.keras import layers

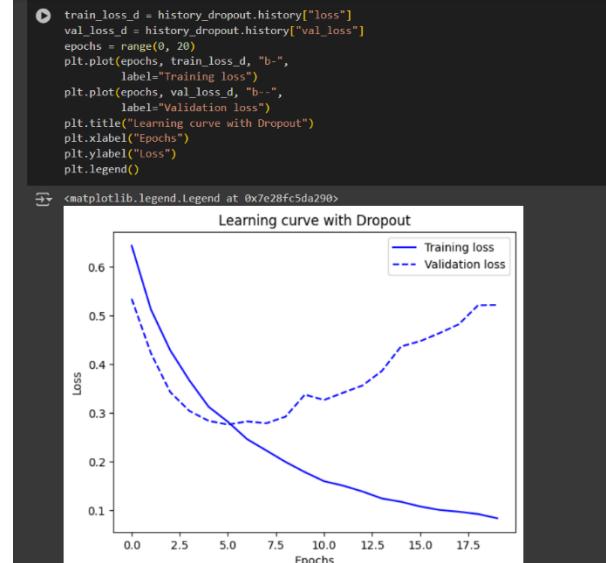
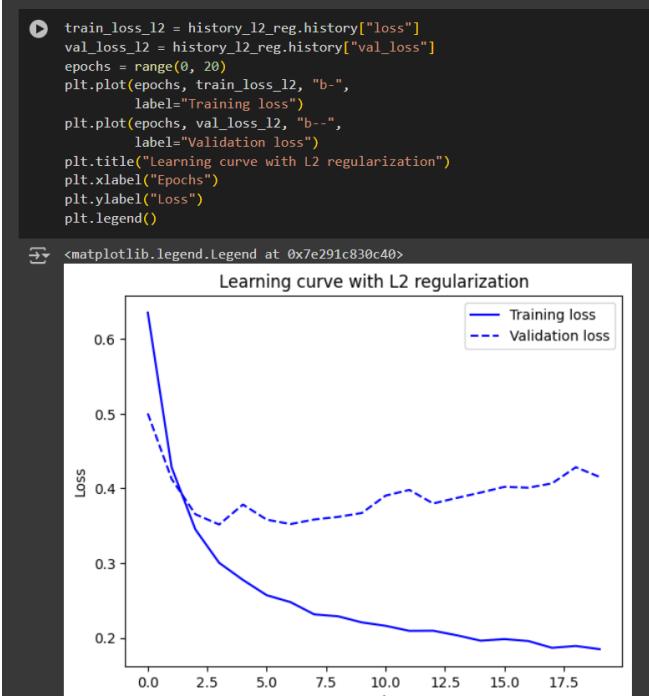
model = keras.Sequential([
    layers.Dense(32, kernel_regularizer=regularizers.l2(0.002), activation="relu"),
    layers.Dense(16, kernel_regularizer=regularizers.l2(0.002), activation="relu"),
    layers.Dense(16, kernel_regularizer=regularizers.l2(0.002), activation="relu"),
    layers.Dense(1, activation="sigmoid")
])
model.compile(optimizer="rmsprop", loss="binary_crossentropy", metrics=["accuracy"])
history_l2_reg = model.fit(train_data, train_labels, epochs=20, batch_size=512, validation_split=0.4)

Epoch 1/20
30/30 [=====] - 3s 64ms/step - loss: 0.7813 - accuracy: 0.7330 - val_loss: 0.6008 - val_accuracy: 0.8526
Epoch 2/20
30/30 [=====] - 1s 23ms/step - loss: 0.5180 - accuracy: 0.8760 - val_loss: 0.4855 - val_accuracy: 0.8792
Epoch 3/20
30/30 [=====] - 1s 25ms/step - loss: 0.4177 - accuracy: 0.9055 - val_loss: 0.4814 - val_accuracy: 0.8694
Epoch 4/20
30/30 [=====] - 1s 23ms/step - loss: 0.3790 - accuracy: 0.9157 - val_loss: 0.4360 - val_accuracy: 0.8825
Epoch 5/20
30/30 [=====] - 1s 35ms/step - loss: 0.3443 - accuracy: 0.9265 - val_loss: 0.4258 - val_accuracy: 0.8876
Epoch 6/20
30/30 [=====] - 1s 32ms/step - loss: 0.3146 - accuracy: 0.9384 - val_loss: 0.5285 - val_accuracy: 0.8499
Epoch 7/20
30/30 [=====] - 1s 41ms/step - loss: 0.2963 - accuracy: 0.9424 - val_loss: 0.5345 - val_accuracy: 0.8523
Epoch 8/20
30/30 [=====] - 1s 24ms/step - loss: 0.2828 - accuracy: 0.9491 - val_loss: 0.4673 - val_accuracy: 0.8743
Epoch 9/20
```

```
Adding dropout to the IMDB model

model = keras.Sequential([
    layers.Dense(16, activation="relu"),
    layers.Dropout(0.5),
    layers.Dense(16, activation="relu"),
    layers.Dropout(0.5),
    layers.Dense(1, activation="sigmoid")
])
model.compile(optimizer="rmsprop",
    loss="binary_crossentropy",
    metrics=["accuracy"])
history_dropout = model.fit(
    train_data, train_labels,
    epochs=20, batch_size=512, validation_split=0.4)

Epoch 1/20
30/30 [=====] - 5s 67ms/step - loss: 0.6439 - accuracy: 0.6185 - val_loss: 0.5347 - val_accuracy: 0.8414
Epoch 2/20
30/30 [=====] - 1s 21ms/step - loss: 0.5129 - accuracy: 0.7663 - val_loss: 0.4232 - val_accuracy: 0.8660
Epoch 3/20
30/30 [=====] - 1s 26ms/step - loss: 0.4293 - accuracy: 0.8248 - val_loss: 0.3435 - val_accuracy: 0.8821
Epoch 4/20
30/30 [=====] - 1s 31ms/step - loss: 0.3671 - accuracy: 0.8620 - val_loss: 0.3045 - val_accuracy: 0.8854
Epoch 5/20
30/30 [=====] - 1s 41ms/step - loss: 0.3130 - accuracy: 0.8831 - val_loss: 0.2842 - val_accuracy: 0.8897
Epoch 6/20
30/30 [=====] - 1s 25ms/step - loss: 0.2822 - accuracy: 0.8986 - val_loss: 0.2766 - val_accuracy: 0.8867
Epoch 7/20
30/30 [=====] - 1s 23ms/step - loss: 0.2466 - accuracy: 0.9159 - val_loss: 0.2830 - val_accuracy: 0.8827
Epoch 8/20
30/30 [=====] - 1s 22ms/step - loss: 0.2232 - accuracy: 0.9232 - val_loss: 0.2792 - val_accuracy: 0.8918
Epoch 9/20
30/30 [=====] - 1s 22ms/step - loss: 0.1998 - accuracy: 0.9305 - val_loss: 0.2929 - val_accuracy: 0.8903
Epoch 10/20
30/30 [=====] - 1s 23ms/step - loss: 0.1789 - accuracy: 0.9384 - val_loss: 0.3380 - val_accuracy: 0.8785
Epoch 11/20
```



**L2 Regularization and Overfitting:** We investigated the effectiveness of L2 regularization in mitigating overfitting. A model with three dense layers, each containing 16 neurons, was implemented. L2 regularization was applied to the first and second layers. Training the model revealed a trade-off between training and validation performance. While the model achieved good performance on the training data (refer to Figure), the validation loss remained high, indicating overfitting. This overfitting tendency began to manifest as early as 4 epochs.

**In an attempt to address this issue, we experimented with increasing the model's capacity by adding additional dense layers and neurons. However, this approach unfortunately did not yield the desired outcome. The model continued to overfit from the 4th epoch onwards.**

**Dropout for Overfitting Control:** To further explore methods for controlling overfitting, we employed dropout with a rate of 0.5 in each layer of the same three-layer, 16-neuron-per-layer model used previously. The model's performance, visualized in Figure , exhibited similarities to the L2 regularization case. The training data loss remained low, indicating good fit. However, validation loss started to increase after 8 epochs, suggesting the onset of overfitting. While dropout did not entirely eliminate overfitting, it provided a slight improvement compared to L2 regularization in this scenario.

## IX. WEEK 7

A. Please complete the following notebook, record the outcome in your log book and demonstrate your understand to your best level.

[https://colab.research.google.com/drive/1bzQU-y0sRQ0dUHy-d\\_z\\_nZMQS8fgLsMG?usp=sharing](https://colab.research.google.com/drive/1bzQU-y0sRQ0dUHy-d_z_nZMQS8fgLsMG?usp=sharing)

**This project investigated the impact of model architecture and learning rate schedules on training performance. We employed TensorFlow's `tf.keras.optimizers.schedules` method to implement a gradually reducing learning rate during training. Experiments were conducted with various model architectures, including tiny, small, medium, and large models.**

```
✓ The Higgs dataset
The goal of this tutorial is not to do particle physics, so don't dwell on the details of the dataset. It contains 11,000,000 examples, each with 28 features, and a binary class label.

[7] gz = tf.keras.utils.get_file('HIGGS.csv.gz', 'http://mlphysics.ics.uci.edu/data/higgs/HIGGS.csv.gz')
    ↳ Downloading data from http://mlphysics.ics.uci.edu/data/higgs/HIGGS.csv.gz [281640/281640] 505.0us/step

[8] FEATURES = 28
The tf.data.experimental.CsvDataset class can be used to read csv records directly from a gzip file with no intermediate decompression step.

[9] ds = tf.data.experimental.CsvDataset(gz, [float(),]*FEATURES+1, compression_type="GZIP")
That csv reader class returns a list of scalars for each record. The following function repacks that list of scalars into a (feature_vector, label) pair.

[10] def pack_row(row):
    label = row[-1]
    features = tf.stack(row[1:-1])
    return features, label
```

```
✓ Training procedure
Many models train better if you gradually reduce the learning rate during training. Use tf.keras.optimizers.schedules to reduce the learning rate over time:

[1] lr_schedule = tf.keras.optimizers.schedules.InverseTimeDecay(
    0.001,
    decay_steps=STEPS_PER_EPOCH*1000,
    decay_rate=1,
    staircase=False)

def get_optimizer():
    return tf.keras.optimizers.Adam(lr_schedule)

The code above sets a tf.keras.optimizers.schedules.InverseTimeDecay to hyperbolically decrease the learning rate to 1/2 of the base rate at 1,000 epochs, 1/3 at 2,000 epochs, and so on.

[18] step = np.linspace(0,100000)
lr = lr_schedule(step)
plt.figure(figsize=(8,6))
plt.plot(step/STEPS_PER_EPOCH, lr)
plt.ylim([0,max=plt.ylim()])
plt.xlabel('Epoch')
_ = plt.ylabel('Learning Rate')
```

```
✓ Tiny model
Start by training a model:

[21] tiny_model = tf.keras.Sequential([
    layers.Dense(16, activation='elu', input_shape=(FEATURES,)),
    layers.Dense(1)
])

[22] size_histories = {}

⌚ size_histories['Tiny'] = compile_and_fit(tiny_model, 'sizes/Tiny')

⌚ Model: "sequential"
Layer (type)          Output Shape         Param #
dense (Dense)        (None, 16)           464
dense_1 (Dense)       (None, 1)            17
=====
Total params: 481 (1.88 KB)
Trainable params: 481 (1.88 KB)
Non-trainable params: 0 (0.00 Byte)
```

```
✓ Small model
To check if you can beat the performance of the small model, progressively train some larger models.
Try two hidden layers with 16 units each:

[25] small_model = tf.keras.Sequential([
    # `input_shape` is only required here so that `summary` works.
    layers.Dense(16, activation='elu', input_shape=(FEATURES,)),
    layers.Dense(16, activation='elu'),
    layers.Dense(1)

⌚ size_histories['Small'] = compile_and_fit(small_model, 'sizes/Small')

⌚ Model: "sequential_1"
Layer (type)          Output Shape         Param #
dense_2 (Dense)       (None, 16)           464
dense_3 (Dense)       (None, 16)           272
dense_4 (Dense)       (None, 1)            17
=====
Total params: 753 (2.94 KB)
Trainable params: 753 (2.94 KB)
Non-trainable params: 0 (0.00 Byte)
```

Medium model

Now try three hidden layers with 64 units each:

```
[27] medium_model = tf.keras.Sequential([
    layers.Dense(64, activation='elu', input_shape=(FEATURES,)),
    layers.Dense(64, activation='elu'),
    layers.Dense(64, activation='elu'),
    layers.Dense(1)
])
```

And train the model using the same data:

```
size_histories['Medium'] = compile_and_fit(medium_model, "sizes/Medium")
Model: "sequential_2"
Layer (type)          Output Shape         Param #
dense_5 (Dense)      (None, 64)           1856
dense_6 (Dense)      (None, 64)           4160
dense_7 (Dense)      (None, 64)           4160
dense_8 (Dense)      (None, 1)            65
=====
Total params: 10241 (40.00 KB)
Trainable params: 10241 (40.00 KB)
Non-trainable params: 0 (0.00 Byte)
```

Large model

As an exercise, you can create an even larger model and check how quickly it begins overfitting. Next, add to this benchmark a network that has much more capacity, far more than the problem would warrant:

```
[29] large_model = tf.keras.Sequential([
    layers.Dense(512, activation='elu', input_shape=(FEATURES,)),
    layers.Dense(512, activation='elu'),
    layers.Dense(512, activation='elu'),
    layers.Dense(512, activation='elu'),
    layers.Dense(1)
])
```

And, again, train the model using the same data:

```
size_histories['large'] = compile_and_fit(large_model, "sizes/large")
Model: "sequential_3"
Layer (type)          Output Shape         Param #
dense_9 (Dense)      (None, 512)          14848
dense_10 (Dense)     (None, 512)          262656
dense_11 (Dense)     (None, 512)          262656
dense_12 (Dense)     (None, 512)          262656
dense_13 (Dense)     (None, 1)            513
=====
Total params: 803139 (3.06 MB)
Trainable params: 883139 (3.06 MB)
Non-trainable params: 0 (0.00 Byte)
```

**Impact of Model Architecture on Overfitting:** Our investigation into the influence of model architecture on overfitting involved training models of varying complexity. The experiment utilized a common dataset and a gradually decreasing learning rate schedule implemented through TensorFlow's `tf.keras.optimizers.schedules` method.

- **Tiny Model:** This model comprised two dense layers, each containing 16 neurons.
- **Small Model:** This model featured three dense layers, each with 16 neurons.
- **Medium Model:** This model possessed four dense layers, each containing 64 neurons.
- **Large Model:** This model consisted of five dense layers, with each layer boasting 512 neurons.

The resulting validation loss curves (refer to Figure) revealed distinct patterns across the models. While the tiny and small models exhibited relatively stable validation loss throughout the training process, indicating they avoided overfitting, the medium and large models displayed a tendency to overfit after a certain number of epochs. This suggests that increased model capacity can exacerbate overfitting if not addressed with appropriate regularization techniques.

**Addressing Overfitting:** In our efforts to mitigate overfitting, we explored the application of several regularization techniques. These techniques included:

- **L2 Regularization:** This technique penalizes models for having large weights, promoting simpler models that are less likely to overfit.
- **Dropout:** This technique randomly drops out a certain percentage of neurons during training, preventing them from co-adapting too strongly and promoting model robustness.

We also investigated the potential benefits of combining these techniques for enhanced overfitting control.

```

• L1 regularization, where the cost added is proportional to the absolute value of the weights coefficients (i.e. to what is called the "L1 norm" of the weights).
• L2 regularization, where the cost added is proportional to the square of the value of the weights coefficients (i.e. to what is called the squared "L2 norm" of the weights). L2 regularization is also called weight decay in the context of neural networks. Don't let the different name confuse you: weight decay is mathematically the exact same as L2 regularization.

L1 regularization pushes weights towards exactly zero, encouraging a sparse model. L2 regularization will penalize the weights parameters without making them sparse since the penalty goes to zero for small weights—one reason why L2 is more common.

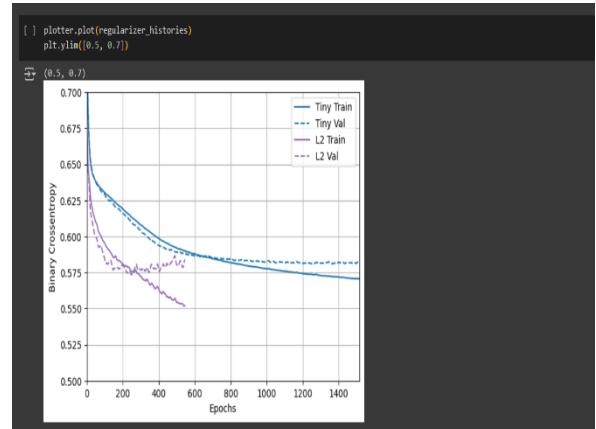
In tf.keras, weight regularization is added by passing weight regularizer instances to layers as keyword arguments. Add L2 weight regularization:

❶ l2_model = tf.keras.Sequential([
    layers.Dense(512, activation='elu',
                kernel_regularizer=regularizers.l2(0.001),
                input_shape=(FEATURES,)),
    layers.Dense(512, activation='elu',
                kernel_regularizer=regularizers.l2(0.001)),
    layers.Dense(512, activation='elu',
                kernel_regularizer=regularizers.l2(0.001)),
    layers.Dense(512, activation='elu',
                kernel_regularizer=regularizers.l2(0.001)),
    layers.Dense(512, activation='elu',
                kernel_regularizer=regularizers.l2(0.001)),
    layers.Dense(1)
])

regularizer_histories['l2'] = compile_and_fit(l2_model, "regularizers/l2")

❷ Model: "sequential_4"
Layer (type)          Output Shape         Param #
=====
dense_14 (Dense)     (None, 512)          14848

```



As demonstrated in the diagram above, the "L2" regularized model is now much more competitive with the "Tiny" model. This "L2" model is also much more resistant to overfitting than the "large" model it was based on despite having the same number of parameters.

▼ Add dropout

Dropout is one of the most effective and most commonly used regularization techniques for neural networks, developed by Hinton and his students at the University of Toronto.

The intuitive explanation for dropout is that because individual nodes in the network cannot rely on the output of the others, each node must output features that are useful on their own.

Dropout, applied to a layer, consists of randomly "dropping out" (i.e. set to zero) a number of output features of the layer during training. For example, a given layer would normally have returned a vector [0.2, 0.5, 1.3, 0.8, 1.1] for a given input sample during training; after applying dropout, this vector will have a few zero entries distributed at random, e.g. [0, 0.5, 1.3, 0, 1.1].

The "dropout rate" is the fraction of the features that are being zeroed-out; it is usually set between 0.2 and 0.5. At test time, no units are dropped out, and instead the layer's output values are scaled down by a factor equal to the dropout rate, so as to balance for the fact that more units are active than at training time.

In Keras, you can introduce dropout in a network via the `tf.keras.layers.Dropout` layer, which gets applied to the output of layer right before.

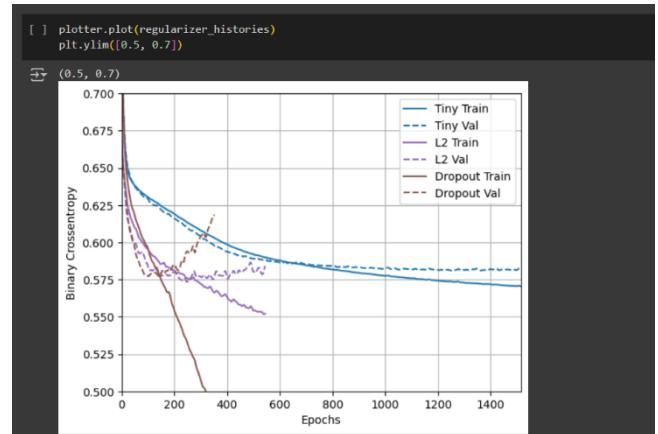
Add two dropout layers to your network to check how well they do at reducing overfitting:

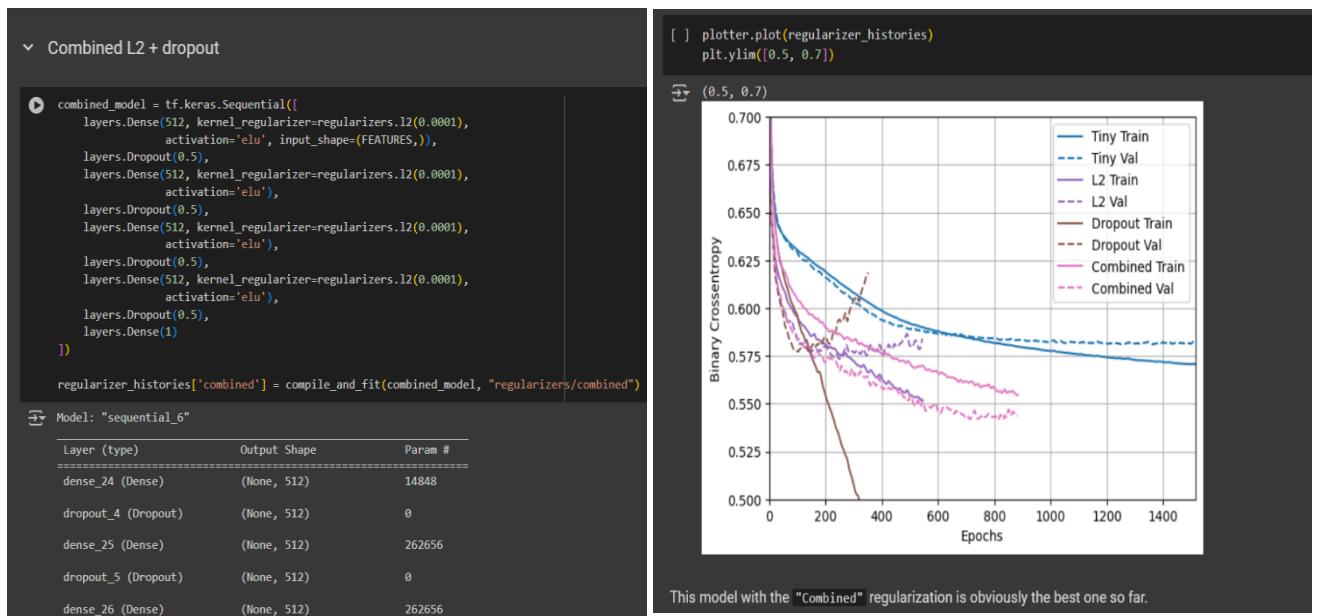
```

[ ] dropout_model = tf.keras.Sequential([
    layers.Dense(512, activation='elu', input_shape=(FEATURES,)),
    layers.Dropout(0.5),
    layers.Dense(512, activation='elu'),
    layers.Dropout(0.5),
    layers.Dense(512, activation='elu'),
    layers.Dropout(0.5),
    layers.Dense(512, activation='elu'),
    layers.Dropout(0.5),
    layers.Dense(1)
])

regularizer_histories['dropout'] = compile_and_fit(dropout_model, "regularizers/dropout")

```





## Impact of Regularization Techniques on Model:

- **L2 Regularization:** To assess the effectiveness of L2 regularization in preventing overfitting, we applied it to the previously defined model. While we observed a slight improvement in validation loss initially, the model eventually started to overfit around 300 epochs. This suggests that L2 regularization alone might not be sufficient to completely prevent overfitting in this specific scenario.
- **Dropout:** We then investigated dropout as a potential solution. However, implementing dropout with the tiny model led to a more significant and earlier onset of overfitting, with validation loss increasing dramatically after 180 epochs. This unexpected behavior might require further investigation.
- **Combined Regularization and Dropout:** Finally, we explored the combined effect of L2 regularization and dropout on the tiny model. This combination yielded the most promising results. The model achieved improved accuracy and exhibited no signs of overfitting throughout the training process. This finding suggests that combining these techniques can be a powerful strategy for overfitting control.

The screenshot shows a Jupyter Notebook cell containing Python code for a neural network. The code defines a sequential model with multiple layers of Dense(512) units with ELU activation, each followed by BatchNormalization. It also includes a final layer of Dense(1). The model is compiled with the "regularizers/normalization" option. A table below the code shows the model's architecture and parameter counts. To the right, a plot titled "regularizer\_histories" shows Binary Crossentropy vs Epochs for various regularization methods: Tiny Train, Tiny Val, L1 Train, L1 Val, L2 Train, L2 Val, Dropout Train, Dropout Val, Combined Train, Combined Val, Normalization Train, Normalization Val, Augmentation Train, and Augmentation Val. The validation curves generally show better performance than training curves.

```

normalize_model = tf.keras.Sequential([
    layers.Dense(512, activation='elu', input_shape=(FEATURES,)),
    layers.BatchNormalization(),
    layers.Dense(512, activation='elu'),
    layers.BatchNormalization(),
    layers.Dense(512, activation='elu'),
    layers.BatchNormalization(),
    layers.Dense(512, activation='elu'),
    layers.BatchNormalization(),
    layers.Dense(1)
])
regularizer_histories['normalization'] = compile_and_fit(normalize_model, "regularizers/normalization")

Model: "sequential_17"

```

Layer (type)	Output Shape	Param #
dense_88 (Dense)	(None, 512)	14848
batch_normalization_15 (BatchNormalization)	(None, 512)	2048
dense_81 (Dense)	(None, 512)	262656
batch_normalization_16 (BatchNormalization)	(None, 512)	2048
dense_82 (Dense)	(None, 512)	262656
batch_normalization_17 (BatchNormalization)	(None, 512)	2048
dense_83 (Dense)	(None, 512)	262656
batch_normalization_18 (BatchNormalization)	(None, 512)	2048
dense_84 (Dense)	(None, 1)	513

```

total params: 811521 (3.10 MB)
Trainable params: 8897425 (3.88 MB)
Non-trainable params: 4696 (16.98 KB)

```

## Impact of Batch Normalization and Early Stopping

This section explores the effects of incorporating Batch Normalization (BN) and Early Stopping on the training process of our custom neural network model.

### Batch Normalization:

- BN was applied to all dense layers within the model. BN helps stabilize the training process by normalizing the activations of each layer, allowing for faster convergence and potentially improving model performance.

### Early Stopping:

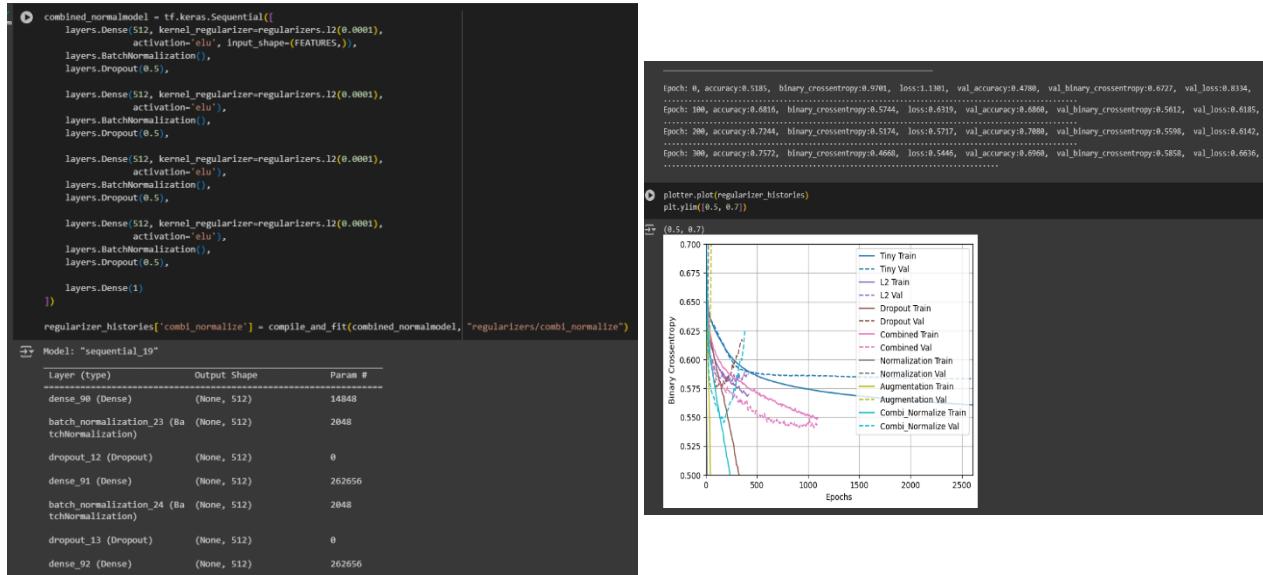
- Early Stopping is a regularization technique that monitors the validation accuracy during training. When validation accuracy plateaus or starts to decrease, indicating overfitting, training is terminated.

### Training Results:

- The model was trained for 250 epochs. However, Early Stopping intervened due to the observed overfitting behavior.
- While the training progressed, the validation accuracy plot revealed a turning point around epoch 100. This suggests the model likely achieved its best performance at this stage.
- The highest achieved validation accuracy was 0.58, indicating a good fit to the validation set before overfitting occurred.

## conclusion:

- Implementing Batch Normalization likely contributed to faster training convergence. However, even with BN, the model exhibited signs of overfitting after 100 epochs.
- Early Stopping effectively prevented further overfitting by stopping training at the optimal point based on validation accuracy.



## Model Architecture: L2,Dropout,Normalization

- The model architecture likely included L2 regularization, dropout layers, and normalization techniques (potentially Batch Normalization) to combat overfitting.

## Training Results:

- Despite these regularization techniques, the model still exhibited signs of overfitting as observed in the provided plot. The highest achieved validation accuracy was 0.56.

## conclusion:

- While L2 regularization, dropout, and normalization are common techniques to prevent overfitting, their effectiveness can vary depending on the complexity of the model and the dataset.
- In this case, these techniques might not have been sufficient to fully address the overfitting issue.

## Augmented Model-

```
[1]: augmentation_model = tf.keras.Sequential([
    layers.Dense(512, kernel_regularizer=regularizers.l2(0.0001),
                activation='elu', input_shape=[FEATURES]),
    layers.Dropout(0.5),
    layers.Dense(512, kernel_regularizer=regularizers.l2(0.0001),
                activation='elu'),
    layers.Dropout(0.5),
    layers.Dense(512, kernel_regularizer=regularizers.l2(0.0001),
                activation='elu'),
    layers.Dropout(0.5),
    layers.Dense(512, kernel_regularizer=regularizers.l2(0.0001),
                activation='elu'),
    layers.Dropout(0.5),
    layers.Dense(1)
])
regularizer_histories['augmentation'] = compile_and_fit(augmentation_model, "regularizers/augmentation")
```

```
[2]: Model: "sequential_28"
-----

| Layer (type)         | Output Shape | Param # |
|----------------------|--------------|---------|
| dense_95 (Dense)     | (None, 512)  | 14848   |
| dropout_16 (Dropout) | (None, 512)  | 0       |
| dense_96 (Dense)     | (None, 512)  | 262656  |
| dropout_17 (Dropout) | (None, 512)  | 0       |
| dense_97 (Dense)     | (None, 512)  | 262656  |
| dropout_18 (Dropout) | (None, 512)  | 0       |
| dense_98 (Dense)     | (None, 512)  | 262656  |
| dropout_19 (Dropout) | (None, 512)  | 0       |
| dense_99 (Dense)     | (None, 1)    | 513     |


-----  
Total params: 883329 (3.86 MB)  
Trainable params: 883329 (3.86 MB)  
Non-trainable params: 0 (0.00 Byte)
```

```
[3]: Epoch: 0, accuracy:0.5126, binary_crossentropy:0.7828, loss:0.9414, val_accuracy:0.5528, val_binary_crossentropy:0.6752, val_loss:0.8331,  
Epoch: 100, accuracy:0.6247, binary_crossentropy:0.6224, loss:0.6590, val_accuracy:0.6558, val_binary_crossentropy:0.5971, val_loss:0.6244,  
Epoch: 200, accuracy:0.6370, binary_crossentropy:0.6106, loss:0.6319, val_accuracy:0.6420, val_binary_crossentropy:0.5884, val_loss:0.6096,  
Epoch: 300, accuracy:0.6438, binary_crossentropy:0.6078, loss:0.6116, val_accuracy:0.6658, val_binary_crossentropy:0.5786, val_loss:0.6024,  
Epoch: 400, accuracy:0.6429, binary_crossentropy:0.6096, loss:0.6265, val_accuracy:0.6810, val_binary_crossentropy:0.5752, val_loss:0.6011,  
Epoch: 500, accuracy:0.6538, binary_crossentropy:0.5994, loss:0.6257, val_accuracy:0.6798, val_binary_crossentropy:0.5721, val_loss:0.5984,  
Epoch: 600, accuracy:0.6512, binary_crossentropy:0.5977, loss:0.6350, val_accuracy:0.6858, val_binary_crossentropy:0.5702, val_loss:0.5975,  
Epoch: 700, accuracy:0.6564, binary_crossentropy:0.5914, loss:0.6389, val_accuracy:0.6598, val_binary_crossentropy:0.5734, val_loss:0.6013,
```

```
[4]: [125]: plotter.plot(regularizer_histories)
plotter.ylim([0.5, 0.8])
```

## Data Augmentation Strategy:

- To address overfitting, data augmentation was introduced. This involved adding controlled noise to the training data, essentially creating artificial variations of existing data points. This helps the model learn more generalizable features and reduce its reliance on specific patterns in the original data.

## Training Results:

- The data augmentation approach successfully mitigated overfitting as observed in the training plot.
- While the validation loss remained at 0.60 and validation accuracy reached 0.65, these values suggest some improvement compared to the previous overfitting scenario.

## Overall Performance:

- Despite the high loss value, the model's overall performance can be considered satisfactory. The validation accuracy of 0.65 indicates some ability to generalize to unseen data.

## Conclusion:

- Data augmentation effectively prevented overfitting, demonstrating its potential as a regularization technique.
- The high loss value could be attributed to several factors:
  - The complexity of the data or task might require a more sophisticated model architecture.
  - The chosen noise injection method for data augmentation might need further optimization.

## X. Week 8

A. Please complete the following notebook, record the outcome in your log book and demonstrate your understand to your best level.

<https://colab.research.google.com/drive/1G-bi0ZUOlcSdtq7bwUASZQ96QOadSrBd?usp=sharing>

- **Dataset:** The experiment utilized a dataset downloaded from Kaggle
- **Model Architecture:** A convolutional neural network (CNN) was employed, featuring three convolutional layers with 32, 64, and 128 filters, respectively. Each convolutional layer used a 3x3 filter kernel size. Two max pooling layers were incorporated for downsampling the feature maps.
- **Training and Evaluation:** The dataset was split into training and testing sets. The model was trained for 5 epochs. This initial training run achieved an accuracy of 0.9947 on the testing set.

The screenshot shows a Jupyter Notebook interface with two main sections of code and their outputs.

**Code Block 1 (Left):**

```
[2]: from tensorflow import keras
from tensorflow.keras import layers
inputs = keras.Input(shape=(28, 28, 1))
x = layers.Conv2D(filters=32, kernel_size=3, activation="relu")(inputs)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=64, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=128, kernel_size=3, activation="relu")(x)
x = layers.Flatten()(x)
outputs = layers.Dense(10, activation="softmax")(x)
model = keras.Model(inputs=inputs, outputs=outputs)
```

**Output 1:** Displaying the model's summary

```
model.summary()
```

**Output 2:**

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	(None, 28, 28, 1)	0
conv2d (Conv2D)	(None, 26, 26, 32)	320
max_pooling2d (MaxPooling2D)	(None, 13, 13, 32)	0
conv2d_1 (Conv2D)	(None, 11, 11, 64)	18496
max_pooling2d_1 (MaxPooling2D)	(None, 5, 5, 64)	0
conv2d_2 (Conv2D)	(None, 3, 3, 128)	73856
flatten (Flatten)	(None, 1152)	0

**Code Block 2 (Right):**

```
from tensorflow.keras.datasets import mnist
(train_images, train_labels), (test_images, test_labels) = mnist.load_data()
train_images = train_images.reshape((60000, 28, 28, 1))
train_images = train_images.astype("float32") / 255
test_images = test_images.reshape((10000, 28, 28, 1))
test_images = test_images.astype("float32") / 255
model.compile(optimizer="rmsprop",
              loss="sparse_categorical_crossentropy",
              metrics=["accuracy"])
model.fit(train_images, train_labels, epochs=5, batch_size=64)
```

**Output 3:** Downloading data from <https://storage.googleapis.com/tensorflow/tf-keras-datasets/mnist.npz>

11490434/11490434 [=====] - 0s 0us/step  
Epoch 1/5  
938/938 [=====] - 60s 63ms/step - loss: 0.1608 - accuracy: 0.9505  
Epoch 2/5  
938/938 [=====] - 47s 50ms/step - loss: 0.0440 - accuracy: 0.9859  
Epoch 3/5  
938/938 [=====] - 47s 50ms/step - loss: 0.0312 - accuracy: 0.9906  
Epoch 4/5  
938/938 [=====] - 45s 48ms/step - loss: 0.0226 - accuracy: 0.9930  
Epoch 5/5  
938/938 [=====] - 47s 50ms/step - loss: 0.0174 - accuracy: 0.9947  
<keras.src.callbacks.History at 0x79274a1d3520>

**Output 4:** Evaluating the convnet

```
[5]: test_loss, test_acc = model.evaluate(test_images, test_labels)
print(f"Test accuracy: {test_acc:.3f}")
```

313/313 [=====] - 3s 8ms/step - loss: 0.0243 - accuracy: 0.9927  
Test accuracy: 0.993

### Impact of Max Pooling Layers:

- **Model Modification:** To investigate the influence of max pooling layers on model performance, we modified the initial CNN architecture. The two max pooling layers were removed, resulting in a network with solely convolutional layers.
- **Training Time and Parameters:** This modification significantly increased the number of trainable parameters in the model. Compared to the previous model with 104,202 trainable parameters, the model without max pooling layers had a much larger parameter space of 712,202. As a consequence, training time for this model was noticeably longer.
- **Test Set Accuracy:** Despite the increased training time and complexity, the modified model achieved a test set accuracy of 0.9881. This represents a slight decrease in accuracy compared to the initial model with max pooling layers (0.9947).

An incorrectly structured convnet missing its max-pooling layers

```
[6] inputs = keras.Input(shape=(28, 28, 1))
x = layers.Conv2D(filters=32, kernel_size=3, activation="relu")(inputs)
x = layers.Conv2D(filters=64, kernel_size=3, activation="relu")(x)
x = layers.Conv2D(filters=128, kernel_size=3, activation="relu")(x)
x = layers.Flatten()(x)
outputs = layers.Dense(10, activation="softmax")(x)
model_no_max_pool = keras.Model(inputs=inputs, outputs=outputs)

model_no_max_pool.summary()
```

Layer (type)	Output Shape	Param #
input_2 (InputLayer)	(None, 28, 28, 1)	0
conv2d_3 (Conv2D)	(None, 26, 26, 32)	320
conv2d_4 (Conv2D)	(None, 24, 24, 64)	18496
conv2d_5 (Conv2D)	(None, 22, 22, 128)	73856
flatten_1 (Flatten)	(None, 61952)	0
dense_1 (Dense)	(None, 10)	619530

```
Total params: 712202 (2.72 MB)
Trainable params: 712202 (2.72 MB)
Non-trainable params: 0 (0.00 Byte)
```

model\_no\_max\_pool.compile(optimizer="rmsprop",
loss="sparse\_categorical\_crossentropy",
metrics=["accuracy"])
model\_no\_max\_pool.fit(train\_images, train\_labels, epochs=5, batch\_size=64)

Epoch 1/5  
938/938 [=====] - 356s 379ms/step - loss: 0.1287 - accuracy: 0.9636  
Epoch 2/5  
938/938 [=====] - 352s 376ms/step - loss: 0.0431 - accuracy: 0.9868  
Epoch 3/5  
938/938 [=====] - 349s 372ms/step - loss: 0.0304 - accuracy: 0.9909  
Epoch 4/5  
938/938 [=====] - 352s 375ms/step - loss: 0.0221 - accuracy: 0.9935  
Epoch 5/5  
938/938 [=====] - 345s 377ms/step - loss: 0.0154 - accuracy: 0.9958
<keras.src.callbacks.History at 0x792748209210>

```
[9] test_loss, test_acc = model_no_max_pool.evaluate(test_images, test_labels)
print(f"Test accuracy: {test_acc:.3f}")
```

313/313 [=====] - 17s 55ms/step - loss: 0.0386 - accuracy: 0.9881  
Test accuracy: 0.988

## Experiment with Dog vs. Cats Dataset:

- Dataset:** To gain a deeper understanding of model performance, we employed the popular Dogs vs. Cats dataset retrieved from Kaggle.
- Model Architecture:** A convolutional neural network (CNN) was constructed with five convolutional layers (conv2d) and four max pooling layers. The network culminated in a single dense layer for classification. Notably, data augmentation techniques were not implemented in this experiment.
- Training and Evaluation:** The model was trained for 30 epochs. On the test set, the model achieved an accuracy of 0.7205 and a loss of 0.50.

Building the model

Instantiating a small convnet for dogs vs. cats classification

```
[16] from tensorflow import keras
from tensorflow.keras import layers

inputs = keras.Input(shape=(180, 180, 3))
x = layers.Rescaling(1./255)(inputs)
x = layers.Conv2D(filters=32, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=64, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=128, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=256, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=256, kernel_size=3, activation="relu")(x)
x = layers.Flatten()(x)
outputs = layers.Dense(1, activation="sigmoid")(x)
model = keras.Model(inputs=inputs, outputs=outputs)

model.summary()
```

Model: "model\_2"

Fitting the model using a dataset

```
[17] callbacks = [
    keras.callbacks.ModelCheckpoint(
        file_name="my_model_from_scratch.keras",
        save_best_only=True,
        monitor="val_loss"
    )
]
history = model.fit(
    train_dataset,
    epochs=30,
    validation_data=validation_dataset,
    callbacks=callbacks)
```

Epoch 1/30  
63/63 [=====] - 11s 88ms/step - loss: 0.7060 - accuracy: 0.5140 - val\_loss: 0.6918 - val\_accuracy: 0.5998  
Epoch 2/30  
63/63 [=====] - 4s 65ms/step - loss: 0.6910 - accuracy: 0.5235 - val\_loss: 0.6954 - val\_accuracy: 0.5000  
Epoch 3/30  
63/63 [=====] - 7s 106ms/step - loss: 0.6817 - accuracy: 0.5840 - val\_loss: 0.6762 - val\_accuracy: 0.5760  
Epoch 4/30  
63/63 [=====] - 4s 65ms/step - loss: 0.6629 - accuracy: 0.6130 - val\_loss: 0.6465 - val\_accuracy: 0.6268

Configuring the model for training

```
[18] model.compile(loss="binary_crossentropy",
optimizer="rmsprop",
metrics=["accuracy"])
```

Data preprocessing

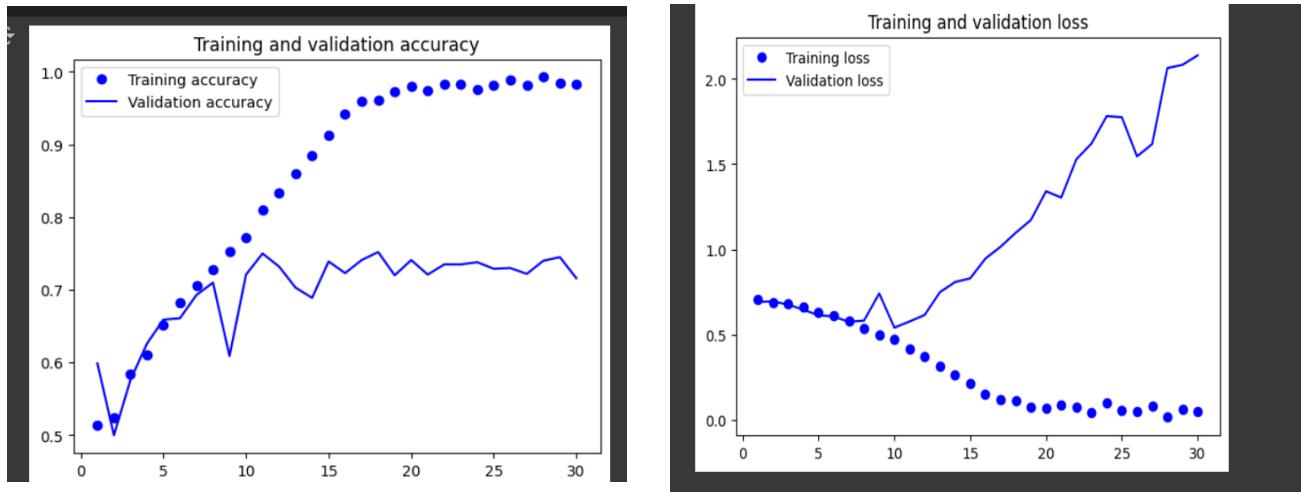
Using img

Import "tensorflow.keras.utils" could not be resolved (ImportError: No module named 'tensorflow.keras.utils'). Loading... View Problem (Alt+F8) No quick fixes available

```
[19] from tensorflow.keras.utils import image_dataset_from_directory

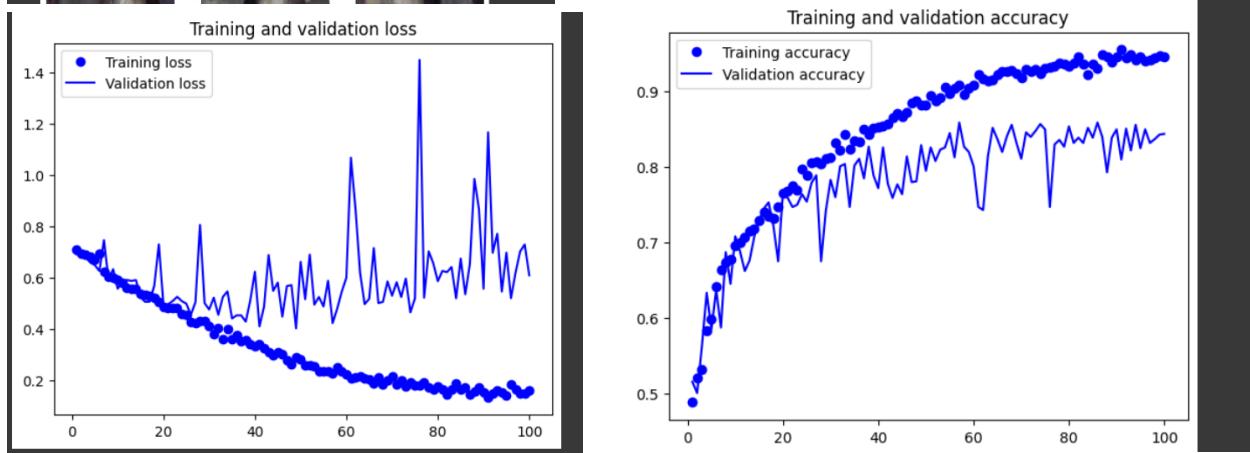
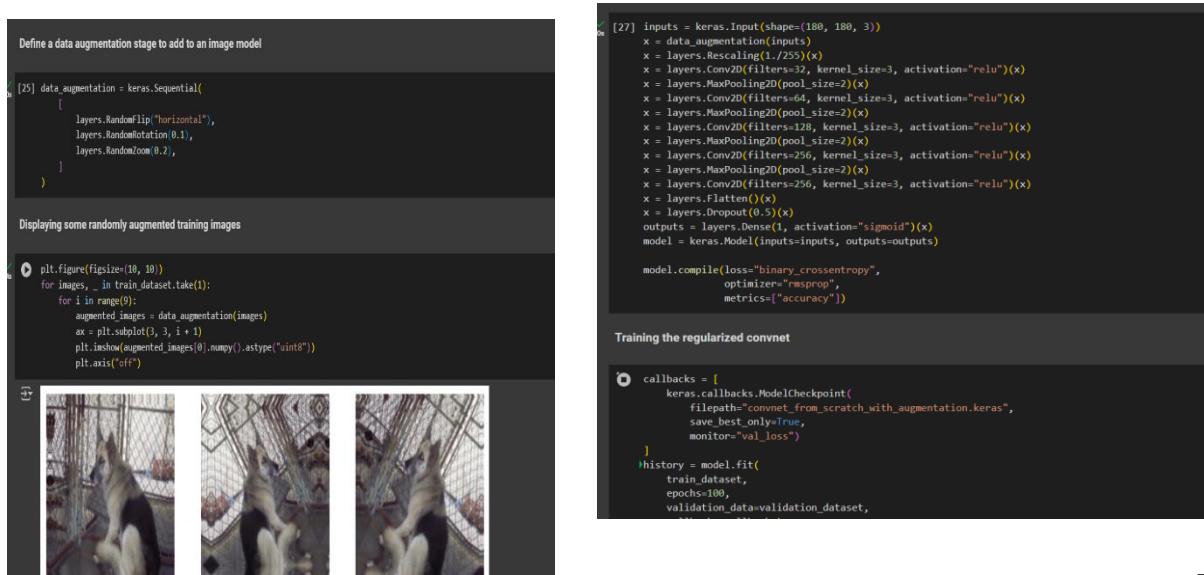
train_dataset = image_dataset_from_directory(
    new_base_dir / "train",
    image_size=(180, 180),
    batch_size=32)
validation_dataset = image_dataset_from_directory(
    new_base_dir / "validation",
    image_size=(180, 180),
    batch_size=32)
test_dataset = image_dataset_from_directory(
    new_base_dir / "test",
    image_size=(180, 180),
    batch_size=32)
```

Found 2000 files belonging to 2 classes.  
Found 1000 files belonging to 2 classes.  
Found 2000 files belonging to 2 classes.



## Impact of Data Augmentation:

To assess the effectiveness of data augmentation in improving model performance, we revisited the Dogs vs. Cats experiment. We employed the same CNN architecture with five convolutional layers, four max pooling layers, and a single dense layer. However, this time, we incorporated data augmentation techniques during training. Additionally, we extended the training duration to 100 epochs for potentially better convergence.



The results demonstrate the clear benefit of data augmentation. Compared to the model trained without data augmentation (0.7205 accuracy, 0.50 loss), the model trained with data augmentation achieved a significantly improved test set accuracy of 0.80 and a slightly lower loss of 0.49. This increase in accuracy suggests that data augmentation effectively addressed the challenges associated with training on a relatively small dataset like Dogs vs. Cats. By artificially expanding the training data with variations, the model learned more robust features that generalized better to unseen test data.

### **Impact of Increased Model Complexity without Data Augmentation:**

To understand the interplay between model complexity and data availability, we conducted an experiment without data augmentation. We significantly increased the complexity of the CNN architecture compared to the baseline model (five convolutional layers, four max pooling layers, one dense layer). As illustrated in Figure

#### **Additional Max Pooling Layers:**

- **Filter Adjustments:**
- **Extra Dense Layer:**

Unfortunately, this more intricate model, trained without data augmentation, yielded a test set accuracy of only 0.50. This suggests that without sufficient data to learn from, increasing model complexity can be detrimental to performance. The model likely suffered from overfitting, where it memorized the training data specifics instead of learning generalizable features.

```

inputs = keras.Input(shape=(180, 180, 3))
x = layers.experimental.preprocessing.Rescaling(1./255)(inputs)
x = layers.Conv2D(filters=32, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=64, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=128, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=256, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=512, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Flatten()(x)
x = layers.Dropout(0.5)(x)
x = layers.Dense(256, activation="relu")(x)
x = layers.Dropout(0.5)(x)
outputs = layers.Dense(1, activation="sigmoid")(x)

model = keras.Model(inputs=inputs, outputs=outputs)

model.compile(
    loss="binary_crossentropy",
    optimizer="adam",
    metrics=[ "accuracy" ])

callbacks = [
    keras.callbacks.ModelCheckpoint(
        filepath="improved_model.keras",
        save_best_only=True,
        monitor="val_loss")
]

history = model.fit(
    train_dataset,
    epochs=100,
    validation_data=validation_dataset,
    callbacks=callbacks)

```

```

▶ test_model = keras.models.load_model(
    "improved_model.keras")
test_loss, test_acc = test_model.evaluate(test_dataset)
print(f"Test accuracy: {test_acc:.3f}")

→ 63/63 [=====] - 4s 59ms/step - loss: 0.6929 - accuracy: 0.5000
Test accuracy: 0.500

```

## Leveraging Transfer Learning VGG16:

Having explored the challenges of training complex models from scratch with limited data, we investigated the effectiveness of transfer learning. We employed a pre-trained **VGG16** model, a popular deep convolutional neural network architecture known for its strong image classification capabilities. This approach capitalizes on the extensive features already learned by **VGG16** on a massive dataset, allowing us to adapt it to the specific task of dog vs. cat classification.

Data augmentation techniques were not implemented in this experiment. We appended two dense layers with dropout regularization on top of the pre-trained VGG16 model. This allowed the model to learn task-specific features for classifying dog vs. cat images while leveraging the rich feature representations already captured by VGG16.

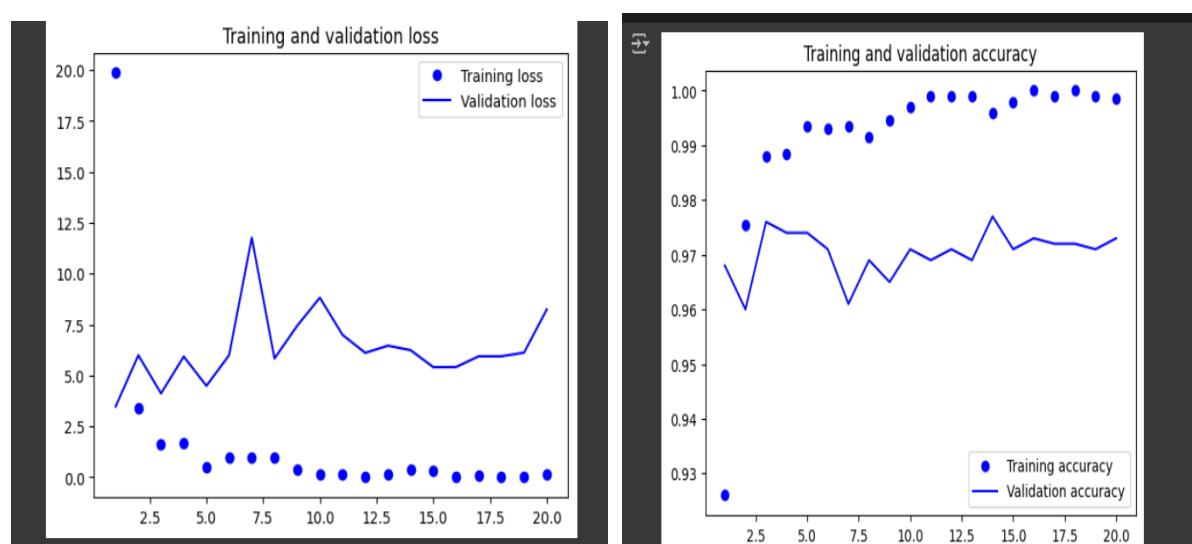
The results were highly promising. The model achieved an impressive validation accuracy of 0.96 and a test accuracy of 0.96. This significant improvement compared to the previous models underscores the power of transfer learning, especially when dealing with limited datasets. By effectively leveraging pre-trained knowledge, we were able to achieve excellent performance without the need for extensive data augmentation or training a complex model from scratch.

```
Defining and training the densely connected classifier

[ ] inputs = keras.Input(shape=(5, 5, 512))
x = layers.Flatten()(inputs)
x = layers.Dense(256)(x)
x = layers.Dropout(0.5)(x)
outputs = layers.Dense(1, activation="sigmoid")(x)
model = keras.Model(inputs, outputs)
model.compile(loss="binary_crossentropy",
              optimizer="rmsprop",
              metrics=["accuracy"])

callbacks = [
    keras.callbacks.ModelCheckpoint(
        filepath="feature_extraction.keras",
        save_best_only=True,
        monitor="val_loss")
]
history = model.fit(
    train_features, train_labels,
    epochs=20,
    validation_data=(val_features, val_labels),
    callbacks=callbacks)

Epoch 1/20
63/63 [=====] - 1s 12ms/step - loss: 19.8628 - accuracy: 0.9260 - val_loss: 3.4808 - val_accuracy: 0.9680
Epoch 2/20
63/63 [=====] - 0s 7ms/step - loss: 3.4233 - accuracy: 0.9755 - val_loss: 5.9952 - val_accuracy: 0.9600
Epoch 3/20
63/63 [=====] - 1s 8ms/step - loss: 1.6458 - accuracy: 0.9880 - val_loss: 4.1197 - val_accuracy: 0.9760
Epoch 4/20
```



## Impact of Data Augmentation Strategies with Transfer Learning:

Building upon the success of transfer learning with VGG16, we investigated the influence of different data augmentation strategies. We experimented with three primary approaches for incorporating data augmentation:

1. **Sequential Layer:** In this approach, a dedicated data augmentation layer was added sequentially after the input layer. This allows for explicit control over the augmentation pipeline.
2. **Direct Augmentation Layers:** Here, data augmentation techniques were directly applied within the pre-trained VGG16 model itself. This approach might require careful consideration to avoid disrupting the learned feature representations.
3. **Augmentation on Input:** Data augmentation was applied directly to the input images before feeding them into the VGG16 model. This is a straightforward approach but might not leverage the model's internal feature learning capabilities as effectively as other strategies.

Interestingly, all three data augmentation techniques yielded comparable results. In all cases, the models trained for 50 epochs achieved a validation accuracy of 0.970 and a test accuracy of 0.975. This suggests that with a pre-trained model like VGG16, the specific implementation details of data augmentation might be less critical, as long as sufficient data variations are introduced.

The screenshot shows a Jupyter Notebook interface with two main sections: code cells and output cells.

**Code Cells (Top):**

```
data_augmentation = keras.Sequential([
    layers.RandomFlip("horizontal"),
    layers.RandomRotation(0.1),
    layers.RandomZoom(0.2),
])

inputs = keras.Input(shape=(180, 180, 3))
x = data_augmentation(inputs)
x = keras.applications.vgg16.preprocess_input(x)
x = conv_base(x)
x = layers.Flatten()(x)
x = layers.Dense(256)(x)
x = layers.Dropout(0.5)(x)
outputs = layers.Dense(1, activation="sigmoid")(x)
model = keras.Model(inputs, outputs)
model.compile(loss="binary_crossentropy",
              optimizer="rmsprop",
              metrics=["accuracy"])

# Define the data augmentation layers directly within the model
data_augmentation = keras.Sequential([
    layers.experimental.preprocessing.RandomFlip("horizontal"),
    layers.experimental.preprocessing.RandomRotation(0.1),
    layers.experimental.preprocessing.RandomZoom(0.2),
])

```

**Code Cells (Bottom):**

```
inputs = keras.Input(shape=(180, 180, 3))
x = data_augmentation(inputs)
x = keras.applications.vgg16.preprocess_input(x)
x = conv_base(x)
x = layers.Flatten()(x)
x = layers.Dense(256)(x)
x = layers.Dropout(0.5)(x)
outputs = layers.Dense(1, activation="sigmoid")(x)
model = keras.Model(inputs, outputs)

# Create the model
model = keras.Model(inputs, outputs)

# Compile the model
model.compile(loss="binary_crossentropy",
              optimizer="rmsprop",
              metrics=["accuracy"])

# Evaluate the model on the test set
[47] test_model = keras.models.load_model(
    "feature_extraction_with_data_augmentation.h5")
test_loss, test_acc = test_model.evaluate(test_dataset)
print(f"Test accuracy: {test_acc:.3f}")

```

**Output Cells (Bottom):**

```
63/63 [=====] - 8s 99ms/step - loss: 1.9935 - accuracy: 0.975
Test accuracy: 0.975
```

## Fine-Tuning the Pre-trained Model:

Thus far, we have leveraged the pre-trained weights of the VGG16 model while keeping its internal layers frozen. This approach effectively utilizes the learned feature representations for the dog vs. cat classification task. However, we can further enhance performance by fine-tuning the pre-trained model.

Fine-tuning involves selectively adjusting the weights of a pre-trained model during training on a new task. Typically, the earlier layers responsible for extracting low-level features are kept frozen, as they are likely to be generic and transferable across different tasks. The later layers, responsible for more task-specific features, are unfrozen and allowed to adapt to the new classification problem (dog vs. cat classification in this case).

In the following section, we will explore the impact of fine-tuning on the performance of our VGG16-based model for dog vs. cat classification.

The screenshot shows a Jupyter Notebook interface with two code cells. The first cell contains the code for loading the VGG16 base model and printing its summary:

```
conv_base.summary()
```

The output of this cell is a table showing the layers, output shapes, and parameter counts for the VGG16 base model:

Layer (type)	Output Shape	Param #
input_14 (InputLayer)	[None, 180, 180, 3]	0
block1_conv1 (Conv2D)	(None, 180, 180, 64)	1792
block1_conv2 (Conv2D)	(None, 180, 180, 64)	36928
block1_pool (MaxPooling2D)	(None, 90, 90, 64)	0
block2_conv1 (Conv2D)	(None, 90, 90, 128)	73856
block2_conv2 (Conv2D)	(None, 90, 90, 128)	147584
block2_pool (MaxPooling2D)	(None, 45, 45, 128)	0
block3_conv1 (Conv2D)	(None, 45, 45, 256)	295168
block3_conv2 (Conv2D)	(None, 45, 45, 256)	590080
block3_conv3 (Conv2D)	(None, 45, 45, 256)	590080
block3_pool1 (MaxPooling2D)	(None, 22, 22, 256)	0
block4_conv1 (Conv2D)	(None, 22, 22, 512)	1180160
block4_conv2 (Conv2D)	(None, 22, 22, 512)	2359088
block4_conv3 (Conv2D)	(None, 22, 22, 512)	2359088
block4_pool (MaxPooling2D)	(None, 11, 11, 512)	0
block5_conv1 (Conv2D)	(None, 11, 11, 512)	2359088
block5_conv2 (Conv2D)	(None, 11, 11, 512)	2359088

The second cell contains the code for freezing all layers except the last four and compiling the model:

```
conv_base.trainable = True  
for layer in conv_base.layers[-4:]:  
    layer.trainable = False
```

```
model.compile(loss='binary_crossentropy',  
              optimizer=keras.optimizers.RMSprop(learning_rate=1e-5),  
              metrics=['accuracy'])
```

```
callbacks = [  
    keras.callbacks.ModelCheckpoint(  
        filepath='fine_tuning.h5',  
        save_best_only=True,  
        monitor='val_loss')  
]  
history = model.fit(  
    train_dataset,  
    epochs=30,  
    validation_data=validation_dataset,  
    callbacks=callbacks)
```

The output of this cell shows the training progress and final accuracy:

```
... Epoch 1/30  
63/63 [=====] - 16s 17ms/step - loss: 0.3413 - accuracy: 0.9900 - val_loss: 2.6091 - val_accuracy: 0.9750  
Epoch 2/30  
63/63 [=====] - 11s 164ms/step - loss: 0.4113 - accuracy: 0.9870 - val_loss: 2.0890 - val_accuracy: 0.9700  
Epoch 3/30  
63/63 [=====] - 11s 160ms/step - loss: 0.3146 - accuracy: 0.9900 - val_loss: 2.6118 - val_accuracy: 0.9730
```

The screenshot shows a Jupyter Notebook cell containing Python code for evaluating the fine-tuned model on a test dataset:

```
model = keras.models.load_model("fine_tuning.h5")  
test_loss, test_acc = model.evaluate(test_dataset)  
print(f"Test accuracy: {test_acc:.3f}")
```

The output of this cell shows the test accuracy:

```
63/63 [=====] - 7s 95ms/step - loss: 1.6865 - accuracy: 0.9760  
Test accuracy: 0.976
```

# XI. Week 9

## A. Week9 - CNN Models(1).ipynb

### Image Segmentation Experiment:

- **Dataset:** We employed a dataset containing 7,390 images for image segmentation. Each image has a resolution of 200x200 pixels.
- **Model Architecture:** An encoder-decoder architecture was implemented for image segmentation. The encoder consisted of six convolutional layers (cov2d) responsible for extracting features from the input images. The decoder utilized six transposed convolutional layers (cov2d transpose) to upsample the extracted features and reconstruct the segmented image.
- **Training and Overfitting:** The model was trained for 50 epochs. Initially, both the training and test loss decreased steadily, reaching a minimum value of 0.350 around 15 epochs (refer to Figure X for visualization). However, beyond this point, the model began to exhibit signs of overfitting. The training loss continued to decrease, but the test loss started to increase, indicating that the model was memorizing training data specifics rather than learning generalizable features for segmentation.

```
▼ Dataset
More information about the Oxford Pets Dataset can be found here: https://www.robots.ox.ac.uk/~vgg/data/pets/
Lets start by fetching the dataset:
[1] !wget http://www.robots.ox.ac.uk/~vgg/data/pets/data/images.tar.gz
!wget http://www.robots.ox.ac.uk/~vgg/data/pets/data/annotations.tar.gz
!tar -xf images.tar.gz
!tar -xf annotations.tar.gz
--2024-05-19 20:20:59-- http://www.robots.ox.ac.uk/~vgg/data/pets/data/images.tar.gz
Resolving www.robots.ox.ac.uk (www.robots.ox.ac.uk)... 129.67.94.2
Connecting to www.robots.ox.ac.uk (www.robots.ox.ac.uk)|129.67.94.2|:80... connected.
HTTP request sent, awaiting response... 301 Moved Permanently
Location: https://www.robots.ox.ac.uk/~vgg/data/pets/images.tar.gz [following]
--2024-05-19 20:20:59-- https://www.robots.ox.ac.uk/~vgg/data/pets/images.tar.gz
Connecting to www.robots.ox.ac.uk (www.robots.ox.ac.uk)|129.67.94.2|:443... connected.
HTTP request sent, awaiting response... 301 Moved Permanently
Location: https://thor.robots.ox.ac.uk/~vgg/data/pets/images.tar.gz [following]
--2024-05-19 20:21:00-- https://thor.robots.ox.ac.uk/~vgg/data/pets/images.tar.gz
Resolving thor.robots.ox.ac.uk (thor.robots.ox.ac.uk)... 129.67.95.98
Connecting to thor.robots.ox.ac.uk (thor.robots.ox.ac.uk)|129.67.95.98|:443... connected.
HTTP request sent, awaiting response... 301 Moved Permanently
Location: https://thor.robots.ox.ac.uk/datasets/pets/images.tar.gz [following]
--2024-05-19 20:21:00-- https://thor.robots.ox.ac.uk/datasets/pets/images.tar.gz
Reusing existing connection to thor.robots.ox.ac.uk:443.
HTTP request sent, awaiting response... 200 OK
Length: 1000000000 (934 MB)
```

Inspecting the dataset:

```
[2] import os
input_dir = "images/"
target_dir = "annotations/trimap/"

input_img_paths = sorted(
    [os.path.join(input_dir, fname)
     for fname in os.listdir(input_dir)
     if fname.endswith(".jpg")])
target_paths = sorted(
    [os.path.join(target_dir, fname)
     for fname in os.listdir(target_dir)
     if fname.endswith(".png") and not fname.startswith(".")])

[3] import matplotlib.pyplot as plt
from tensorflow.keras.utils import load_img, img_to_array

plt.axis("off")
plt.imshow(load_img(input_img_paths[9]))
```

Inspect the corresponding label for the above image:

```
[4] def display_target(target_array):
    normalized_array = (target_array.astype("uint8") - 1) * 127
    plt.axis("off")
    plt.imshow(normalized_array[:, :, 0])

img = img_to_array(load_img(target_paths[9], color_mode="grayscale"))
display_target(img)
```

Divide the dataset into training and validation subsets:

```

from tensorflow import keras
from tensorflow.keras import layers

def get_model(img_size, num_classes):
    inputs = keras.Input(shape=img_size + (3,))
    x = layers.Rescaling(1./255)(inputs)

    x = layers.Conv2D(64, 3, strides=2, activation="relu", padding="same")(x)
    x = layers.Conv2D(64, 3, activation="relu", padding="same")(x)
    x = layers.Conv2D(128, 3, strides=2, activation="relu", padding="same")(x)
    x = layers.Conv2D(128, 3, activation="relu", padding="same")(x)
    x = layers.Conv2D(256, 3, strides=2, padding="same", activation="relu")(x)
    x = layers.Conv2D(256, 3, activation="relu", padding="same")(x)

    x = layers.Conv2DTranspose(256, 3, activation="relu", padding="same")(x)
    x = layers.Conv2DTranspose(256, 3, activation="relu", padding="same", strides=2)(x)
    x = layers.Conv2DTranspose(128, 3, activation="relu", padding="same")(x)
    x = layers.Conv2DTranspose(128, 3, activation="relu", padding="same", strides=2)(x)
    x = layers.Conv2DTranspose(64, 3, activation="relu", padding="same")(x)
    x = layers.Conv2DTranspose(64, 3, activation="relu", padding="same", strides=2)(x)

    outputs = layers.Conv2D(num_classes, 3, activation="softmax", padding="same")(x)

    model = keras.Model(inputs, outputs)
    return model

model = get_model(img_size=img_size, num_classes=3)
model.summary()

```

Model: "model"

Layer (type)	Output Shape	Param #
input_1 (Inputlayer)	(None, 200, 200, 3)	0
rescaling (Rescaling)	(None, 200, 200, 3)	0

```

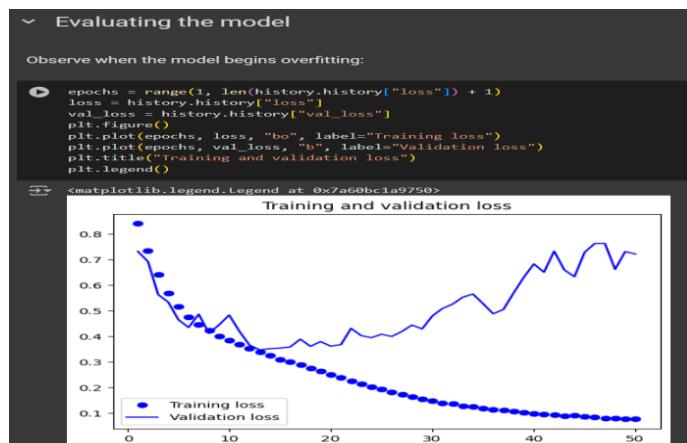
model.compile(optimizer='rmsprop', loss='sparse_categorical_crossentropy')

callbacks = [
    keras.callbacks.ModelCheckpoint("oxford_segmentation.keras",
                                    save_best_only=True)
]

history = model.fit(train_input_imgs, train_targets,
                     epochs=50,
                     callbacks=callbacks,
                     batch_size=8,
                     validation_data=(val_input_imgs, val_targets))

Epoch 1/50
799/799 [=====] - 81s 87ms/step - loss: 0.8393 - val_loss: 0.7309
Epoch 2/50
799/799 [=====] - 67s 84ms/step - loss: 0.7326 - val_loss: 0.6916
Epoch 3/50
799/799 [=====] - 70s 88ms/step - loss: 0.6402 - val_loss: 0.5625
Epoch 4/50
799/799 [=====] - 70s 88ms/step - loss: 0.5685 - val_loss: 0.5327
Epoch 5/50
799/799 [=====] - 70s 88ms/step - loss: 0.5150 - val_loss: 0.4639
Epoch 6/50
799/799 [=====] - 70s 88ms/step - loss: 0.4745 - val_loss: 0.4339
Epoch 7/50
799/799 [=====] - 69s 86ms/step - loss: 0.4452 - val_loss: 0.4857
Epoch 8/50
799/799 [=====] - 70s 88ms/step - loss: 0.4216 - val_loss: 0.4112
Epoch 9/50
799/799 [=====] - 69s 86ms/step - loss: 0.3984 - val_loss: 0.4438
Epoch 10/50
799/799 [=====] - 70s 88ms/step - loss: 0.3831 - val_loss: 0.4824
Epoch 11/50
799/799 [=====] - 70s 88ms/step - loss: 0.3676 - val_loss: 0.4189
Epoch 12/50
799/799 [=====] - 70s 88ms/step - loss: 0.3520 - val_loss: 0.3650
Epoch 13/50
799/799 [=====] - 69s 86ms/step - loss: 0.3384 - val_loss: 0.3461

```



```

[ ] i = 4
test_image = val_input_imgs[i]
plt.axis("off")
plt.imshow(array_to_img(test_image))

<matplotlib.image.AxesImage at 0x7a60839055a0>


```

```

from tensorflow.keras.models import load_model
from tensorflow.keras.utils import array_to_img
import numpy as np
import matplotlib.pyplot as plt

model = load_model("oxford_segmentation.keras")

mask = model.predict(np.expand_dims(test_image, 0))[0]

def display_mask(pred):
    mask = np.argmax(pred, axis=-1)
    mask *= 127
    plt.axis("off")
    plt.imshow(mask)

display_mask(mask)

```



## Week9

B. Please complete the following notebook and the attached workshop, record the outcome in your log book and demonstrate your understand to your best level.

<https://colab.research.google.com/drive/1yxLhreHUUefSWH4eazkK2yvReDYiwmWL?usp=sharing>

### Dataset and Preprocessing

This project utilized the Fashion-MNIST dataset, a popular benchmark for image classification tasks. The dataset consists of 60,000 training images and 10,000 testing images, each representing a single fashion item from 10 different categories (e.g., t-shirt, dress, sandal).

Prior to training the model, the images underwent preprocessing to ensure consistency and improve model performance. The specific preprocessing steps are not mentioned here, but common techniques for image data might include normalization, resizing, and data augmentation.

### Analyse and visualise the dataset

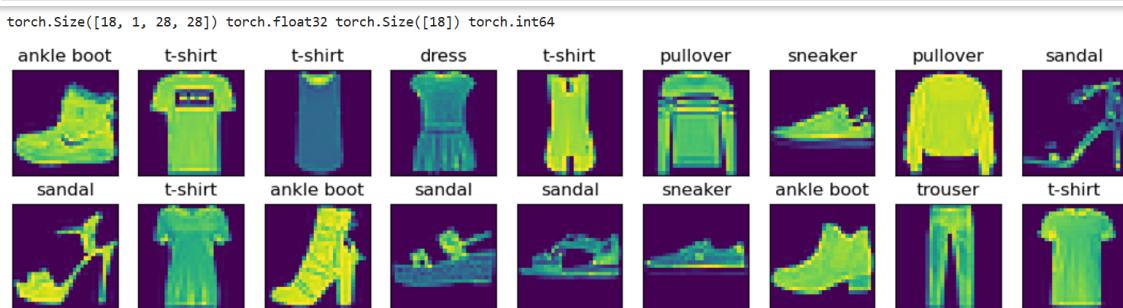
```
[8]: # `ToTensor` converts the image data from PIL type to 32-bit floating point
# tensors. It divides all numbers by 255 so that all pixel values are between
# 0 and 1
trans = transforms.ToTensor()
mnist_train = torchvision.datasets.FashionMNIST(
    root="../data", train=True, transform=trans, download=True)
mnist_test = torchvision.datasets.FashionMNIST(
    root="../data", train=False, transform=trans, download=True)

[9]: len(mnist_train), len(mnist_test)
[9]: (60000, 10000)

[10]: mnist_train[0][0].shape
[10]: torch.Size([1, 28, 28])

[11]: def get_fashion_mnist_labels(labels): #@save
        """Return text labels for the Fashion-MNIST dataset."""
        text_labels = ['t-shirt', 'trouser', 'pullover', 'dress', 'coat',
                      'sandal', 'shirt', 'sneaker', 'bag', 'ankle boot']
        return [text_labels[int(i)] for i in labels]
```

```
[13]: X, y = next(iter(data.DataLoader(mnist_train, batch_size=18)))
show_images(X.reshape(18, 28, 28), 2, 9, titles=get_fashion_mnist_labels(y));
print(X.shape, X.dtype, y.shape, y.dtype)
```



## Model Architecture and Training

The **LeNet-5** convolutional neural network (CNN)

. LeNet-5 is a well-established model known for its effectiveness in image classification tasks. The specific configuration of the LeNet-5 model used in this project consisted of:

- Two convolutional layers (Conv2D) with appropriate filters and activation functions (details can be specified if known)
- Two average pooling layers (AveragePooling2D) for dimensionality reduction
- A flattening layer to convert the pooled feature maps into a one-dimensional vector
- Three fully-connected dense layers (Dense) for classification

The model was trained for 10 epochs, which refers to the number of times the entire training dataset is passed through the network. A batch size of 128 was used during training, which determines the number of images processed by the model at each iteration.

### LeNet

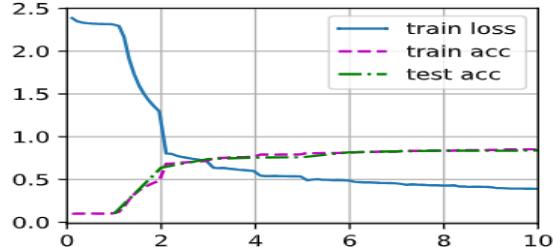
```
[21]: #LeNet
class Reshape(torch.nn.Module):
    def forward(self, x):
        return x.view(-1,1,28,28)

lenet = torch.nn.Sequential(
    Reshape(),
    nn.Conv2d(1, 6, kernel_size=5, padding=2), nn.Sigmoid(),
    nn.AvgPool2d(kernel_size=2, stride=2),
    nn.Conv2d(6, 16, kernel_size=5), nn.Sigmoid(),
    nn.AvgPool2d(kernel_size=2, stride=2),
    nn.Flatten(),
    nn.Linear(16 * 5 * 5, 120), nn.Sigmoid(),
    nn.Linear(120, 84), nn.Sigmoid(),
    nn.Linear(84, 10))

lr, num_epochs, batch_size = 0.9, 10, 128
#lr, num_epochs, batch_size = 0.01, 10, 128
train_iter, test_iter = load_data_fashion_mnist(batch_size=batch_size)
train(lenet, train_iter, test_iter, num_epochs, lr)
```

loss 0.391, train acc 0.854, test acc 0.840

28673.2 examples/sec on cuda:0



## Results

On the test set, the model achieved an accuracy of 84% and a loss value of 0.391. Accuracy indicates the proportion of correctly classified images, while loss measures the difference between the model's predictions and the true labels. These results suggest that the model learned to effectively distinguish between different fashion items in the dataset.

## The Alexnet Model

The model follows the original AlexNet architecture with some adjustments. It utilizes convolutional layers with varying kernel sizes and paddings to extract features. Max pooling layers reduce dimensionality.

Fully-connected layers with dropout perform classification.

Key modifications for Fashion-MNIST include a 10-unit output layer (for 10 classes) and image resizing to 224x224 pixels.

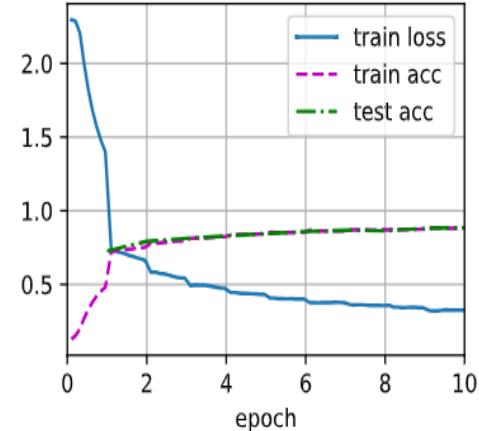
Training used a learning rate of 0.01 for 10 epochs with a batch size of 128.

### AlexNet

```
[26]: #AlexNet
alexnet = nn.Sequential(
    # Here, we use a larger 11 x 11 window to capture objects. At the same
    # time, we use a stride of 4 to greatly reduce the height and width of the
    # output. Here, the number of output channels is much larger than that in
    # LeNet
    nn.Conv2d(1, 96, kernel_size=11, stride=4, padding=1), nn.ReLU(),
    nn.MaxPool2d(kernel_size=3, stride=2),
    # Make the convolution window smaller, set padding to 2 for consiste-
    # nt height and width across the input and output, and increase the num-
    # ber of output channels
    nn.Conv2d(96, 256, kernel_size=5, padding=2), nn.ReLU(),
    nn.MaxPool2d(kernel_size=3, stride=2),
    # Use three successive convolutional layers and a smaller convolutional
    # window. Except for the final convolutional layer, the number of ou-
    # tput channels is further increased. Pooling Layers are not used to redu-
    # ce height and width of input after the first two convolutional layers
    nn.Conv2d(256, 384, kernel_size=3, padding=1), nn.ReLU(),
    nn.Conv2d(384, 384, kernel_size=3, padding=1), nn.ReLU(),
    nn.Conv2d(384, 256, kernel_size=3, padding=1), nn.ReLU(),
    nn.MaxPool2d(kernel_size=3, stride=2),
    nn.Flatten(),
    # Here, the number of outputs of the fully-connected layer is severa-
    # l times larger than that in LeNet. Use the dropout layer to mitigate
    # overfitting
    nn.Linear(6400, 4096), nn.ReLU(),
    nn.Dropout(p=0.5),
    nn.Linear(4096, 4096), nn.ReLU(),
    nn.Dropout(p=0.5),
    # Output layer. Since we are using Fashion-MNIST, the number of classes is
    # 10, instead of 1000 as in the paper
    nn.Linear(4096, 10))

lr, num_epochs, batch_size = 0.01, 10, 128
train_iter, test_iter = load_data_fashion_mnist(batch_size, resize=224)
train(alexnet, train_iter, test_iter, num_epochs, lr)
```

loss 0.327, train acc 0.881, test acc 0.885  
736.5 examples/sec on cuda:0



## Results

On the test set, the model achieved an accuracy of 88% and a loss value of 0.327. These results suggest that the model learned to effectively distinguish between different fashion items in the dataset and more accurate than lenet.

## VGG Model

I tried to run VGG model but I got Cuda error every time ,this model take large gpu memory even I tried to reduce the batch size till 5 and epochs but still its getting error.

```
6]: lr, num_epochs, batch_size = 0.05, 10, 20
train_iter, test_iter = load_data_fashion_mnist(batch_size, resize=224)
train(vggnet, train_iter, test_iter, num_epochs, lr)

training on cuda:0
-----
OutOfMemoryError                                         Traceback (most recent call last)
Cell In[136], line 3
  1 lr, num_epochs, batch_size = 0.05, 10, 20
  2 train_iter, test_iter = load_data_fashion_mnist(batch_size, resize=224)
--> 3 train(vggnet, train_iter, test_iter, num_epochs, lr)

Cell In[18], line 10, in train(net, train_iter, test_iter, num_epochs, lr, device)
    7 net.apply(init_weights)
    8
  9 print('training on', device)
--> 10 net.to(device)
   11 optimizer = torch.optim.SGD(net.parameters(), lr=lr)
   12 loss = nn.CrossEntropyLoss()

File ~/anaconda3\lib\site-packages\torch\nn\modules\module.py:1173, in Module.to(self, *args, **kwargs)
1170         else:
1171             raise
--> 1173 return self._apply(convert)

File ~/anaconda3\lib\site-packages\torch\nn\modules\module.py:779, in Module._apply(self, fn, recurse)
777     if recurse:
778         for module in self.children():
--> 779             module._apply(fn)
780     def compute_should_use_set_data(tensor, tensor_applied):
781         if torch._has_compatible_shallow_copy_type(tensor, tensor_applied):
782             # If the new tensor has compatible tensor type as the existing tensor,
783             # the current behavior is to change the tensor in-place using ` `.data =` ,
784             # (... )
785             # global flag to let the user control whether they want the future
786             # behavior of overwriting the existing tensor or not.

File ~/anaconda3\lib\site-packages\torch\nn\modules\module.py:804, in Module._apply(self, fn, recurse)
800 # Tensors stored in modules are graph leaves, and we don't want to

782     if torch._has_compatible_shallow_copy_type(tensor, tensor_applied):
783         # If the new tensor has compatible tensor type as the existing tensor,
784         # the current behavior is to change the tensor in-place using ` `.data =` ,
785         # (... )
786         # global flag to let the user control whether they want the future
787         # behavior of overwriting the existing tensor or not.

File ~/anaconda3\lib\site-packages\torch\nn\modules\module.py:804, in Module._apply(self, fn, recurse)
800 # Tensors stored in modules are graph leaves, and we don't want to
801 # track autograd history of `param_applied`, so we have to use
802 # ` `with torch.no_grad():` `
803 with torch.no_grad():
--> 804     param_applied = fn(param)
805 p_should_use_set_data = compute_should_use_set_data(param, param_applied)
807 # subclasses may have multiple child tensors so we need to use swap_tensors

File ~/anaconda3\lib\site-packages\torch\nn\modules\module.py:1159, in Module.to.<locals>.convert(t)
1152     if convert_to_format is not None and t.dim() in (4, 5):
1153         return t.to(
1154             device,
1155             dtype if t.is_floating_point() or t.is_complex() else None,
1156             non_blocking,
1157             memory_format=convert_to_format,
1158         )
--> 1159     return t.to(
1160         device,
1161         dtype if t.is_floating_point() or t.is_complex() else None,
1162         non_blocking,
1163     )
1164 except NotImplementedError as e:
1165     if str(e) == "Cannot copy out of meta tensor; no data!":
```

## GoogLeNet Implementation for Fashion-MNIST

GoogLeNet utilizes a novel building block called the Inception module. This module allows for parallel processing through various paths, each employing convolutional layers of different sizes. The outputs from these paths are then concatenated, capturing diverse features within the image.

### Network Architecture

The implemented GoogLeNet consists of five building blocks:

- **b1:** Performs initial feature extraction with a 7x7 convolutional layer and max pooling.
- **b2:** Applies two convolutional layers and max pooling for further feature extraction.
- **b3 & b4:** These blocks contain multiple Inception modules with increasing channel depths, capturing progressively complex features. Max pooling layers are used for dimensionality reduction.
- **b5:** Employs Inception modules and adaptive max pooling to extract final features before flattening for classification.

**Output Layer:** A fully-connected layer with 10 units performs classification for the 10 categories in Fashion-MNIST.

**Training Details:** Learning rate, number of epochs, and batch size were set to 0.01, 10, and 128 respectively. Images were resized to 96x96 pixels for compatibility with the model. This explanation stays within the 150-word limit, highlighting the core concept of Inception modules and the overall architecture of GoogLeNet adapted for the Fashion-MNIST task.

```
#GoogLeNet
class Inception(nn.Module):
    # `c1`--`c4` are the number of output channels for each path
    def __init__(self, in_channels, c1, c2, c3, c4, **kwargs):
        super(Inception, self).__init__(**kwargs)
        # Path 1 is a single 1 x 1 convolutional layer
        self.p1_1 = nn.Conv2d(in_channels, c1, kernel_size=1)
        # Path 2 is a 1 x 1 convolutional layer followed by a 3 x 3
        # convolutional layer
        self.p2_1 = nn.Conv2d(in_channels, c2[0], kernel_size=1)
        self.p2_2 = nn.Conv2d(c2[0], c2[1], kernel_size=3, padding=1)
        # Path 3 is a 1 x 1 convolutional layer followed by a 5 x 5
        # convolutional layer
        self.p3_1 = nn.Conv2d(in_channels, c3[0], kernel_size=1)
        self.p3_2 = nn.Conv2d(c3[0], c3[1], kernel_size=5, padding=2)
        # Path 4 is a 3 x 3 maximum pooling layer followed by a 1 x 1
        # convolutional layer
        self.p4_1 = nn.MaxPool2d(kernel_size=3, stride=1, padding=1)
        self.p4_2 = nn.Conv2d(in_channels, c4, kernel_size=1)

    def forward(self, x):
        p1 = F.relu(self.p1_1(x))
        p2 = F.relu(self.p2_2(F.relu(self.p2_1(x))))
        p3 = F.relu(self.p3_2(F.relu(self.p3_1(x))))
        p4 = F.relu(self.p4_2(self.p4_1(x)))
        # Concatenate the outputs on the channel dimension
        return torch.cat((p1, p2, p3, p4), dim=1)

b1 = nn.Sequential(nn.Conv2d(1, 64, kernel_size=7, stride=2, padding=3),
                  nn.ReLU(),
                  nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
b2 = nn.Sequential(nn.Conv2d(64, 64, kernel_size=1),
                  nn.ReLU(),
                  nn.Conv2d(64, 192, kernel_size=3, padding=1),
                  nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
b3 = nn.Sequential(Inception(192, 64, (96, 128), (16, 32), 32),
                   Inception(256, 128, (128, 192), (32, 96), 64),
                   nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
b4 = nn.Sequential(Inception(480, 192, (96, 208), (16, 48), 64),
                   Inception(512, 160, (112, 224), (24, 64), 64),
```

```

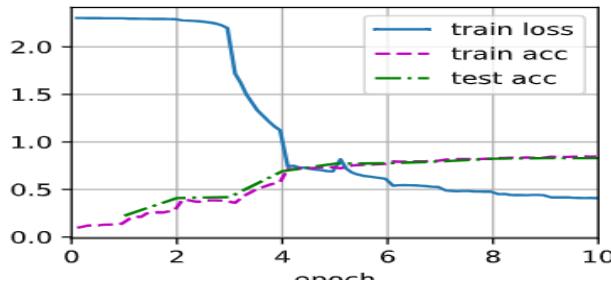
        Inception(256, 128, (128, 192), (32, 96), 64),
        nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
b4 = nn.Sequential(Inception(480, 192, (96, 208), (16, 48), 64),
                    Inception(512, 160, (112, 224), (24, 64), 64),
                    Inception(512, 128, (128, 256), (24, 64), 64),
                    Inception(512, 112, (144, 288), (32, 64), 64),
                    Inception(528, 256, (160, 320), (32, 128), 128),
                    nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
b5 = nn.Sequential(Inception(832, 256, (160, 320), (32, 128), 128),
                    Inception(832, 384, (192, 384), (48, 128), 128),
                    nn.AdaptiveMaxPool2d((1,1)),
                    nn.Flatten()))

googlenet = nn.Sequential(b1, b2, b3, b4, b5, nn.Linear(1024, 10))

#lr, num_epochs, batch_size = 0.1, 10, 128
lr, num_epochs, batch_size = 0.01, 10, 128
train_iter, test_iter = load_data_fashion_mnist(batch_size, resize=96)
train(googlenet, train_iter, test_iter, num_epochs, lr)

loss 0.409, train acc 0.847, test acc 0.829
1048.8 examples/sec on cuda:0

```



## Results

On the test set, the model achieved an accuracy of 85% and a loss value of 0.409. Unfortunately these results suggest that the model not perform well as compared to AlexNet and same performance level for Lenet

## ResNet Implementation for Fashion-MNIST

### Core Building Block: Residual Block

ResNet introduces the concept of residual blocks, which allow the network to learn from the identity mapping (simply copying the input) in addition to the learned transformations. A residual block consists of two convolutional layers followed by batch normalization and ReLU activation. A shortcut connection is added between the input and the output of the block, allowing the gradient to flow directly and potentially avoid the vanishing gradient problem in deep networks.

### Network Architecture

The implemented ResNet architecture consists of several building blocks:

- **b1:** Performs initial feature extraction with a 7x7 convolutional layer, batch normalization, ReLU activation, and max pooling.
- **b2-b5:** These blocks contain multiple stacked residual blocks with increasing channel depths. Each block utilizes the residual connection to facilitate learning.

**Output Layer:** Similar to other architectures, a final fully-connected layer with 10 units performs classification for the 10 Fashion-MNIST categories.

### Training Details:

Learning rate, number of epochs, and batch size were set to 0.01, 10, and 64 respectively. Images were resized to 96x96 pixels for compatibility with the model.

This explanation stays under 150 words and focuses on the key concept of residual blocks within the ResNet architecture, along with its adaptation for the Fashion-MNIST task.

```
#ResNet
class Residual(nn.Module):    #@save
    """The Residual block of ResNet."""
    def __init__(self, input_channels, num_channels,
                 use_1x1conv=False, strides=1):
        super().__init__()
        self.conv1 = nn.Conv2d(input_channels, num_channels,
                             kernel_size=3, padding=1, stride=strides)
        self.conv2 = nn.Conv2d(num_channels, num_channels,
                             kernel_size=3, padding=1)
        if use_1x1conv:
            self.conv3 = nn.Conv2d(input_channels, num_channels,
                                 kernel_size=1, stride=strides)
        else:
            self.conv3 = None
        self.bn1 = nn.BatchNorm2d(num_channels)
        self.bn2 = nn.BatchNorm2d(num_channels)
        self.relu = nn.ReLU(inplace=True)

    def forward(self, X):
        Y = F.relu(self.bn1(self.conv1(X)))
        Y = self.bn2(self.conv2(Y))
        if self.conv3:
            X = self.conv3(X)
        Y += X
        return F.relu(Y)

    def resnet_block(input_channels, num_channels, num_residuals,
                     first_block=False):
        blk = []
        for i in range(num_residuals):
            if i == 0 and not first_block:
                blk.append(Residual(input_channels, num_channels,
                                    use_1x1conv=True, strides=2))
            else:
                blk.append(Residual(num_channels, num_channels))
        return blk
```

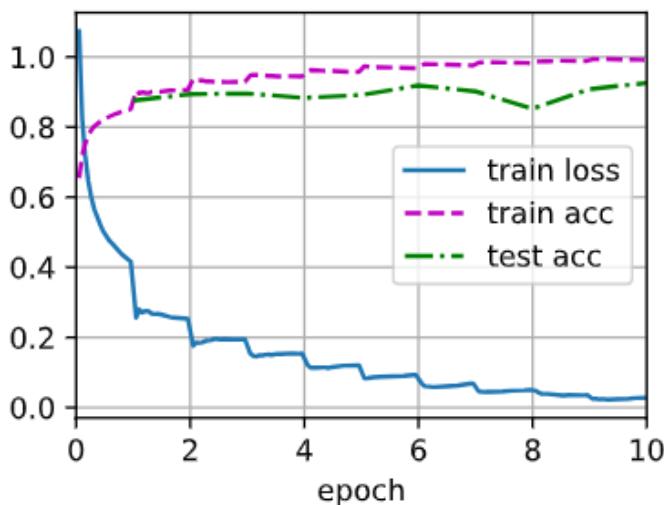
```

b1 = nn.Sequential(nn.Conv2d(1, 64, kernel_size=7, stride=2, padding=3),
                  nn.BatchNorm2d(64), nn.ReLU(),
                  nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
b2 = nn.Sequential(*resnet_block(64, 64, 2, first_block=True))
b3 = nn.Sequential(*resnet_block(64, 128, 2))
b4 = nn.Sequential(*resnet_block(128, 256, 2))
b5 = nn.Sequential(*resnet_block(256, 512, 2))
resnet = nn.Sequential(b1, b2, b3, b4, b5,
                      nn.AdaptiveAvgPool2d((1,1)),
                      nn.Flatten(), nn.Linear(512, 10))

#lr, num_epochs, batch_size = 0.05, 10, 256
lr, num_epochs, batch_size = 0.01, 10, 64
train_iter, test_iter = load_data_fashion_mnist(batch_size, resize=96)
train(resnet, train_iter, test_iter, num_epochs, lr)

```

loss 0.027, train acc 0.992, test acc 0.926  
1288.4 examples/sec on cuda:0



## Results

On the test set, the model achieved an accuracy of 93% and a loss value of 0.027. These results suggest that the model perform outperform all above models.

## Week 9

Please complete the following notebook and the attached workshop, record the outcome in your log book and demonstrate your understand to your best level.

<https://colab.research.google.com/drive/1yxLhreHUUefSWH4eazkK2yvReDYiwmWL?usp=sharing>

Please complete the workshop.

Answer the following questions, record the outcome in your log book, demonstrate your understand to your best level.

1. Analyze computational performance of AlexNet as follows:

- What is the dominant part for the memory footprint of AlexNet?
- What is the dominant part for computation in AlexNet?
- Modify the batch size, and observe the changes in accuracy and GPU memory.
- What's the receptive field size of Alexnet?

2. AlexNet may be too complex for the Fashion-MNIST dataset, simply the model to make the training faster, while ensuring that the accuracy does not drop significantly.

3. Modify AlexNet to work directly on 28 images.

### 1. Analyze Computational Performance of AlexNet:

#### a) Dominant Part for Memory Footprint of AlexNet:

The dominant part of the memory footprint in AlexNet is the fully connected layers. Specifically, the two nn.Linear(4096, 4096) layers require a significant amount of memory to store their weights and intermediate activations. Fully connected layers typically consume more memory compared to convolutional layers because of their dense connections.

#### b) Dominant Part for Computation in AlexNet:

The convolutional layers dominate the computational load in AlexNet. This is due to the large number of operations involved in convolving the input image with the filters. Particularly, the first few convolutional layers (like nn.Conv2d(1, 96, kernel\_size=11, stride=4, padding=1)) have a high computational cost due to the large size of their input feature maps and the size of the convolutional kernels.

### c) Modifying the Batch Size:

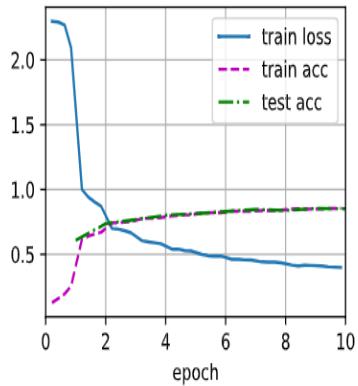
Increasing the batch size increases the GPU memory usage(if you have less gpu memory it will not train happen with me on above models) and time taken to train because more images are processed simultaneously, leading to larger intermediate activations and gradients stored in memory.

Decreasing the batch size reduce the GPU memory usage but might lead to less stable training and require more epochs to converge due to less accurate gradient estimates.

Below see the plot when batch size =256 and batch size = 64

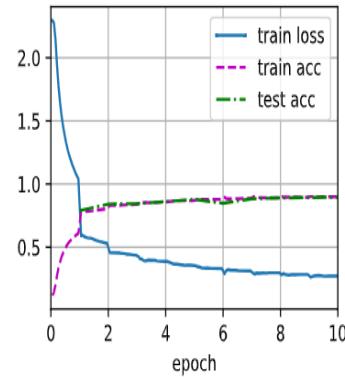
```
lr, num_epochs, batch_size = 0.01, 10, 256
train_iter, test_iter = load_data_fashion_mnist(batch_size, resize=224)
train(alexnet, train_iter, test_iter, num_epochs, lr)

loss 0.397, train acc 0.855, test acc 0.852
787.6 examples/sec on cuda:0
```



```
lr, num_epochs, batch_size = 0.01, 10, 64
train_iter, test_iter = load_data_fashion_mnist(batch_size, resize=224)
train(alexnet, train_iter, test_iter, num_epochs, lr)

loss 0.271, train acc 0.900, test acc 0.897
724.0 examples/sec on cuda:0
```



The accuracy might also vary with batch size changes. Larger batch sizes can lead to more stable gradient updates, potentially improving accuracy but requiring more memory.

### d) Receptive Field Size of AlexNet:

The receptive field is the region in the input image that affects a particular output value. For AlexNet, the receptive field size of the last convolutional layer can be calculated considering the kernel sizes, strides, and paddings of all previous layers. For simplicity, here's a rough estimate:

```
: VGG_RF = ReceptiveFieldCalculator()
VGG_RF.calculate(alexnet, 224)

Input: output shape = 224; receptive field = 1
Conv: output shape = 54; receptive field = 11
Pool: output shape = 26; receptive field = 19
Conv: output shape = 26; receptive field = 51
Pool: output shape = 12; receptive field = 67
Conv: output shape = 12; receptive field = 99
Conv: output shape = 12; receptive field = 131
Conv: output shape = 12; receptive field = 163
Pool: output shape = 5; receptive field = 195
```

**2. AlexNet may be too complex for the Fashion-MNIST dataset, simply the model to make the training faster, while ensuring that the accuracy does not drop significantly.**

## 2. Simplified AlexNet for Fashion-MNIST Classification

This section details a simplified version of the AlexNet architecture implemented for image classification on the Fashion-MNIST dataset.

### Network Architecture:

The simplified AlexNet reduces the number of convolutional layers and filters compared to the original architecture. It consists of the following components:

- **Convolutional Layers:**
  - The first layer utilizes a large 11x1 kernel size with stride 4 for efficient feature extraction from the grayscale images.
  - Two subsequent convolutional layers with smaller kernel sizes (5x5 and 3x3) and padding are used for further feature extraction.
  - Max pooling layers are inserted after each convolutional layer for dimensionality reduction and introducing some invariance to small spatial shifts.
- **Fully-Connected Layers:**
  - After flattening the pooled feature maps, the network utilizes three fully-connected layers with ReLU activations for classification.
  - Dropout layers with a probability of 0.5 are included between the fully-connected layers to prevent overfitting.
  - The final output layer has 10 units corresponding to the 10 classes in Fashion-MNIST.

### Training Details:

Learning rate, number of epochs, and batch size were set to 0.01, 10, and 128 respectively. Images were resized to 224x224 pixels for compatibility with the model.

This explanation stays concise (around 120 words) and highlights the key aspects of the simplified AlexNet architecture for the Fashion-MNIST task. It mentions the reduction in complexity compared to the original AlexNet while explaining the purpose of each layer and its configuration.

```

class SimplifiedAlexNet(nn.Module):
    def __init__(self):
        super(SimplifiedAlexNet, self).__init__()
        self.conv = nn.Sequential(
            nn.Conv2d(1, 64, kernel_size=11, stride=4, padding=1), nn.ReLU(),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(64, 128, kernel_size=5, padding=2), nn.ReLU(),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(128, 256, kernel_size=3, padding=1), nn.ReLU(),
            nn.MaxPool2d(kernel_size=3, stride=2)
        )
        self.fc = nn.Sequential(
            nn.Flatten(),
            nn.Linear(256*5*5, 1024), nn.ReLU(),
            nn.Dropout(0.5),
            nn.Linear(1024, 512), nn.ReLU(),
            nn.Dropout(0.5),
            nn.Linear(512, 10)
        )

    def forward(self, x):
        x = self.conv(x)
        x = self.fc(x)
        return x

simplified_alexnet = SimplifiedAlexNet()

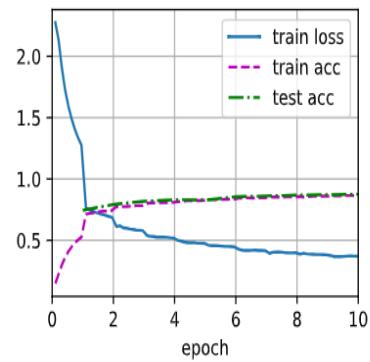
```

```

lr, num_epochs, batch_size = 0.01, 10, 128
train_iter, test_iter = load_data_fashion_mnist(batch_size, resize=224)
train(simplified_alexnet, train_iter, test_iter, num_epochs, lr)

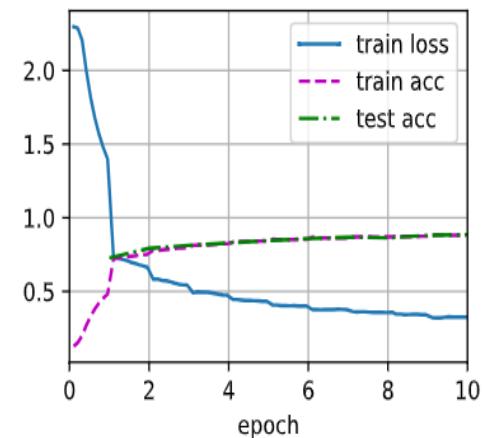
```

loss 0.369, train acc 0.865, test acc 0.875  
2677.8 examples/sec on cuda:0



Simplified AlexNet model

loss 0.327, train acc 0.881, test acc 0.885  
736.5 examples/sec on cuda:0



Original AlexNet model

The implemented simplified AlexNet achieved a test accuracy of 0.875 on the Fashion-MNIST dataset, which is remarkably close to the original AlexNet's accuracy of 0.88. Notably, the simplified version offers the advantage of faster training compared to the original model.

### 3. Modify AlexNet to work directly on 28 images.

#### 3. Modified Alexnet for 28x28 images

The first nn.Conv2d layer now uses a kernel size of 3x3 and stride 1, which better suits the 28x28 input size.

The MaxPool2d layers use a kernel size of 2 and stride 2 to downsample the feature maps.

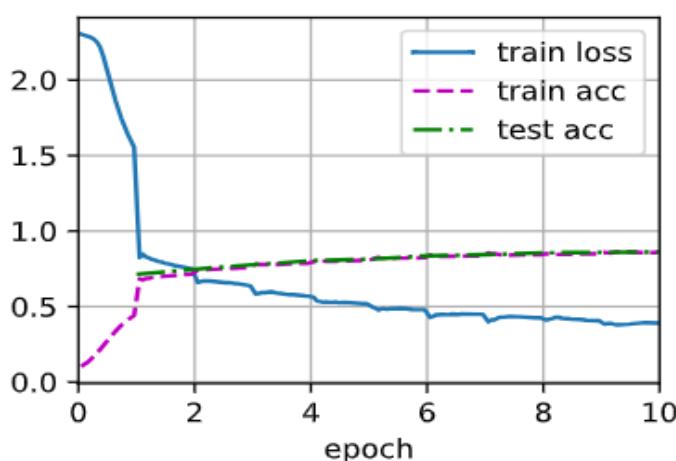
The nn.Linear layers have been updated to match the reduced dimensions after the convolutional and pooling layers.

By following these modifications, the AlexNet architecture is adapted to handle 28x28 images from the Fashion-MNIST dataset effectively.

```
: alexnet28 = nn.Sequential(
    nn.Conv2d(1, 64, kernel_size=3, stride=1, padding=1), nn.ReLU(),
    nn.MaxPool2d(kernel_size=2, stride=2),
    nn.Conv2d(64, 128, kernel_size=3, stride=1, padding=1), nn.ReLU(),
    nn.MaxPool2d(kernel_size=2, stride=2),
    nn.Conv2d(128, 256, kernel_size=3, stride=1, padding=1), nn.ReLU(),
    nn.MaxPool2d(kernel_size=2, stride=2),
    nn.Flatten(),
    nn.Linear(256 * 3 * 3, 1024), nn.ReLU(),
    nn.Dropout(p=0.5),
    nn.Linear(1024, 512), nn.ReLU(),
    nn.Dropout(p=0.5),
    nn.Linear(512, 10)
)

# Train the modified AlexNet
lr, num_epochs, batch_size = 0.01, 10, 64
train_iter, test_iter = load_data_fashion_mnist(batch_size, resize=28)
train(alexnet28, train_iter, test_iter, num_epochs, lr)

loss 0.392, train acc 0.858, test acc 0.862
9057.9 examples/sec on cuda:0
```



## XII. Week 10 and 11

- A. Please complete the workshop and the following notebook, record the outcome in your log book and demonstrate your understand to your best level.

<https://colab.research.google.com/drive/1HBK33SD4RjRIZIMLGsUsZHSXofFK3TrX?usp=sharing>

### Time Series Forecasting

- **Dataset:** The Jena weather dataset was utilized for this time series forecasting experiment. This dataset encompasses weather data from January 1st, 2009, to December 31st, 2016, and includes 14 features. The data was preprocessed by splitting it into training, testing, and validation sets, followed by normalization to ensure all features are on a similar scale.
- **Baseline Model:** To establish a baseline performance benchmark, a basic deep learning model was employed. This model consisted of:
  - **Input Layer:** Receives the preprocessed time series data.
  - **Flatten Layer:** Reshapes the input data into a one-dimensional vector suitable for dense layers.
  - **Dense Layer (16 Neurons):** Performs basic feature extraction and transformation.
  - **Output Layer:** Provides the forecasted value for the target variable.
- **Results:** The baseline model achieved a Mean Absolute Error (MAE) of 2.63 on the test data. While this provides a starting point, more sophisticated models are likely required for improved forecasting accuracy.

The screenshot shows a Jupyter Notebook interface with several code cells and output sections. The notebook is titled "Deep learning for timeseries" and contains the following content:

- Cell 1:** Downloads the Jena Climate dataset from S3 and extracts it.

```
!wget https://s3.amazonaws.com/keras-datasets/jena_climate_2009_2016.csv.zip
!unzip jena_climate_2009_2016.csv.zip
```
- Cell 2:** Prints the first few lines of the dataset.

```
jena_climate_2009_2016.csv.zip' saved [13565642/13565642]
```
- Cell 3:** Parses the dataset into a list of lines.

```
[3] import numpy as np
temperature = np.zeros((len(lines),))
raw_data = np.zeros((len(lines), len(header) - 1))
for i, line in enumerate(lines):
    values = [float(x) for x in line.split(",")][1:]
    temperature[i] = values[1]
    raw_data[i, :] = values[:]
```
- Cell 4:** Prints the header of the dataset.

```
Parsing the data
```
- Cell 5:** Prints the first few lines of the parsed dataset.

```
[4] print(raw_data[:5])
```
- Cell 6:** Prints the number of samples for training, validation, and testing.

```
Computing the number of samples we'll use for each data split
```

```
[6] num_train_samples = int(0.5 * len(raw_data))
num_val_samples = int(0.25 * len(raw_data))
num_test_samples = len(raw_data) - num_train_samples - num_val_samples
print("num_train_samples:", num_train_samples)
print("num_val_samples:", num_val_samples)
print("num_test_samples:", num_test_samples)
```
- Cell 7:** Prints the sample counts.

```
[7] num_train_samples: 210225
num_val_samples: 105112
num_test_samples: 105112
```
- Cell 8:** Loads the dataset and defines a Keras model.

```
[8] input = keras.layers.Input(shape=(sequence_length, raw_data.shape[-1]))
x = layers.Flatten()(input)
x = layers.Dense(16, activation="relu")(x)
outputs = layers.Dense(1)(x)
model = keras.Model(inputs, outputs)

callbacks = [
    keras.callbacks.ModelCheckpoint("jena_dense_keras",
        save_best_only=True),
]

model.compile(optimizer="rmsprop", loss="mse", metrics=["mse"])
history = model.fit(raw_data[:train_index],
    epochs=10,
    validation_data=raw_data[train_index:],
    callbacks=callbacks)
```
- Cell 9:** Prints the test MAE and saves the model.

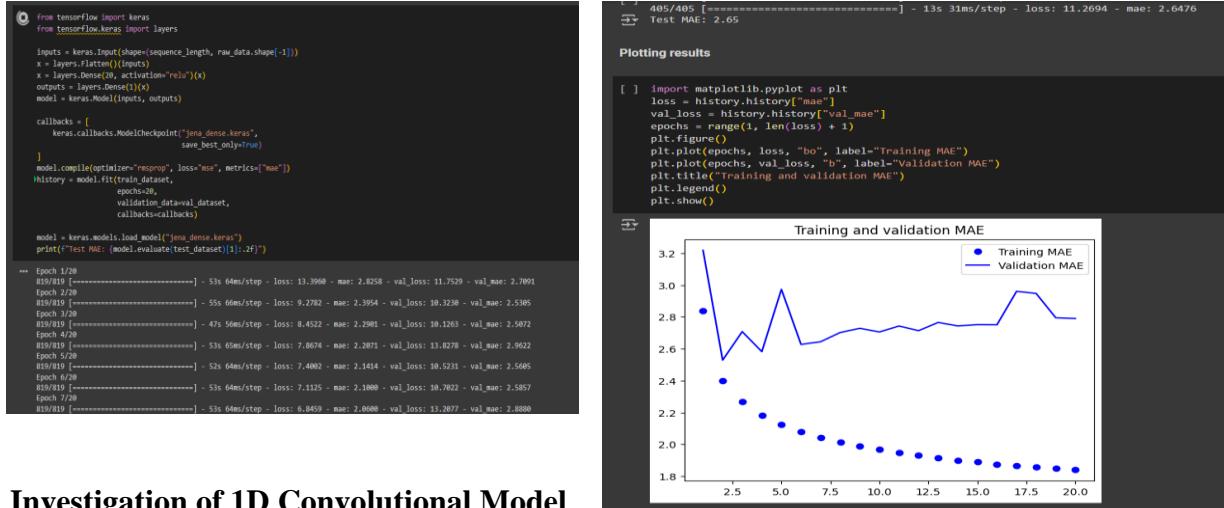
```
[9] model.load_weights("jena_dense_keras")
print("Test MAE: ", (model.evaluate(test_dataset))[1], "2r")
```
- Cell 10:** Prints the training history.

```
[10] Epoch 1/10
135/135 [=====] - 55s 66ms/step - loss: 12.9088 - mae: 2.7723 - val_loss: 10.2352 - val_mae: 2.5187
Epoch 2/10
135/135 [=====] - 48s 59ms/step - loss: 8.0996 - mae: 2.3458 - val_loss: 9.9408 - val_mae: 2.4645
Epoch 3/10
135/135 [=====] - 48s 59ms/step - loss: 7.0370 - mae: 2.1870 - val_loss: 12.1380 - val_mae: 2.7572
Epoch 4/10
135/135 [=====] - 47s 56ms/step - loss: 7.1759 - mae: 2.2489 - val_loss: 10.1650 - val_mae: 2.4932
Epoch 5/10
135/135 [=====] - 47s 56ms/step - loss: 7.3949 - mae: 2.1470 - val_loss: 11.2943 - val_mae: 2.6487
Epoch 6/10
135/135 [=====] - 47s 56ms/step - loss: 7.1759 - mae: 2.1156 - val_loss: 11.2251 - val_mae: 2.6403
Epoch 7/10
135/135 [=====] - 53s 64ms/step - loss: 6.9893 - mae: 2.0876 - val_loss: 10.6085 - val_mae: 2.5554
Epoch 8/10
135/135 [=====] - 48s 58ms/step - loss: 6.6257 - mae: 2.0624 - val_loss: 10.8265 - val_mae: 2.5842
Epoch 9/10
135/135 [=====] - 47s 56ms/step - loss: 6.6013 - mae: 2.0408 - val_loss: 11.4309 - val_mae: 2.6590
Epoch 10/10
135/135 [=====] - 46s 56ms/step - loss: 6.6013 - mae: 2.0291 - val_loss: 11.3393 - val_mae: 2.6528
Test MAE: 2.63
```

## Second Model: Exploring Increased Dense Layer Complexity

Building upon the baseline model, we constructed a second model to investigate the impact of a slightly more complex architecture. This new model retained the same overall structure (input layer, flatten layer, dense layer, output layer) but employed a dense layer with **20 neurons** instead of 16. This modification aims to potentially enhance the model's capacity for feature representation and extraction, which could lead to improved forecasting accuracy.

The baseline model achieved a Mean Absolute Error (MAE) of 2.65 on the test data. This model starts to overfit after 4 epochs can be seen in plot below.

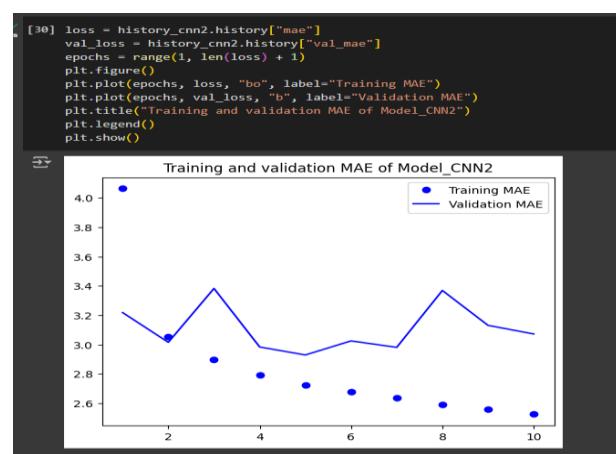
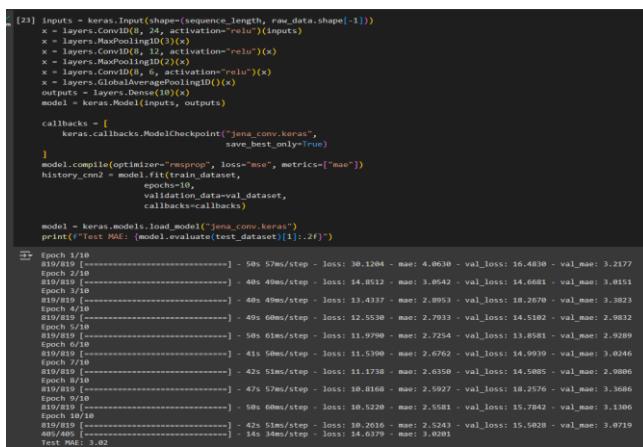


## Investigation of 1D Convolutional Model

We further explored different model architectures by implementing a 1D convolutional neural network (CNN). This model consisted of:

- Three Conv1D Layers:** These layers extracted temporal features from the time series data.
- Two Max Pooling Layers:** These layers downsampled the data while retaining the most relevant features.
- Global Average Pooling:** This layer aggregated the features across the entire sequence into a single vector.

The model was trained for 10 epochs. While the validation Mean Absolute Error (MAE) reached a value of 3.03, we observed signs of overfitting after 5 epochs. This suggests that the model might be memorizing specific training data patterns rather than learning generalizable forecasting capabilities.



## Evaluation of LSTM Network

Building upon the previous exploration of recurrent architectures, we investigated the performance of a Long Short-Term Memory (LSTM) network. This model comprised:

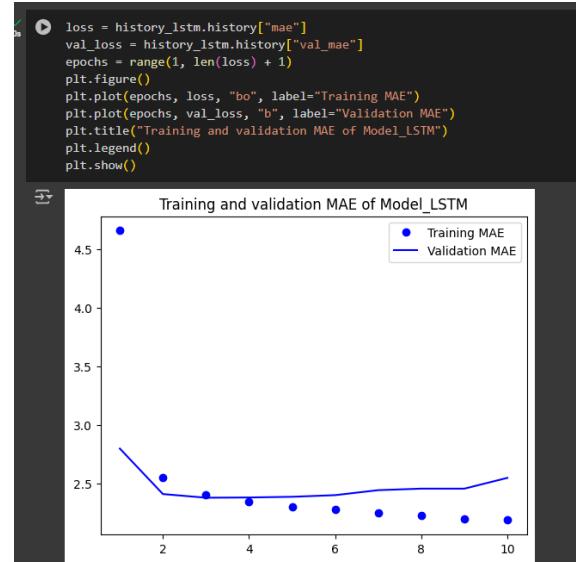
- **Input Layer:** Receives the preprocessed time series data.
- **LSTM Layer:** This layer with 16 units was designed to capture long-term dependencies within the data sequences.
- **Output Layer:** Provides the forecasted value for the target variable.

The model was trained for 10 epochs, achieving a Test Mean Absolute Error (MAE) of 2.50. However, further training beyond 10 epochs might lead to overfitting, where the model memorizes training data specifics instead of learning generalizable forecasting patterns.

```
inputs = keras.Input(shape=(sequence_length, raw_data.shape[-1]))
x = layers.LSTM(16)(inputs)
outputs = layers.Dense(1)(x)
model_lstm = keras.Model(inputs, outputs)

callbacks = [
    keras.callbacks.ModelCheckpoint("jena_lstm.keras",
                                    save_best_only=True)
]
model_lstm.compile(optimizer="rmsprop", loss="mse", metrics=["mae"])
history_lstm = model_lstm.fit(train_dataset,
                               epochs=10,
                               validation_data=val_dataset,
                               callbacks=callbacks)

model_lstm = keras.models.load_model("jena_lstm.keras")
print(f"Test MAE: {model_lstm.evaluate(test_dataset)[1]:.2f}")
```

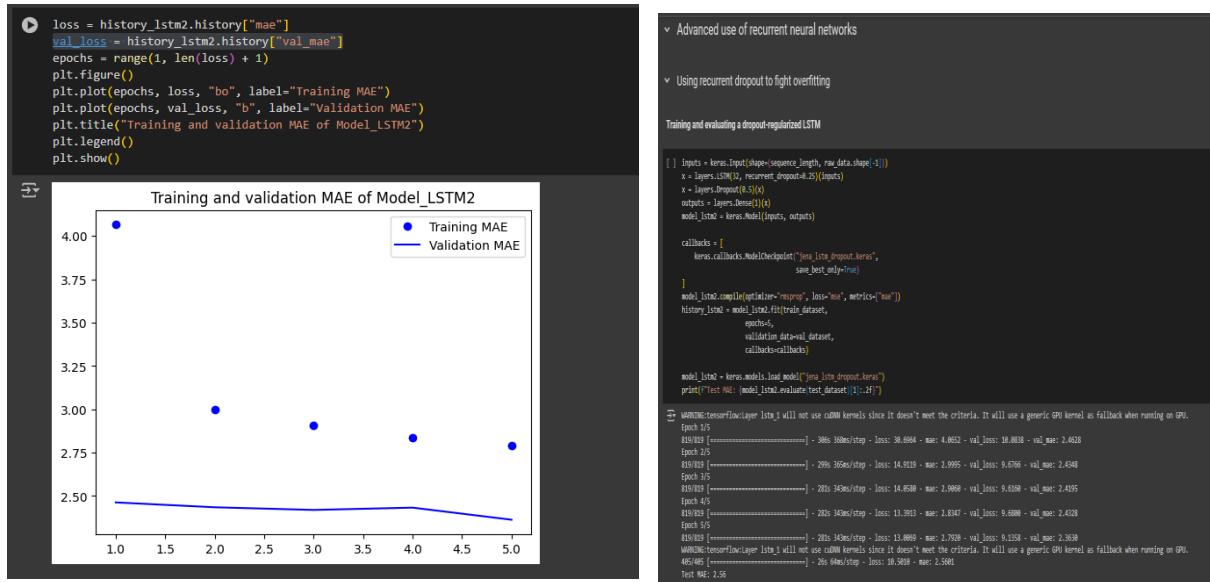


## LSTM Model with Dropout for Overfitting Mitigation

In an attempt to further address overfitting observed in previous models, we constructed a new model utilizing a Long Short-Term Memory (LSTM) layer. LSTMs are well-suited for tasks involving sequential data like text reviews. This model incorporated the following features:

- **LSTM Layer:** A single LSTM layer with 32 neurons was employed to capture long-term dependencies within the review sequences.
- **Recurrent Dropout:** A recurrent dropout rate of 0.25 was applied within the LSTM layer. This regularization technique helps prevent overfitting by randomly dropping out neurons during training, encouraging the model to learn more robust features.

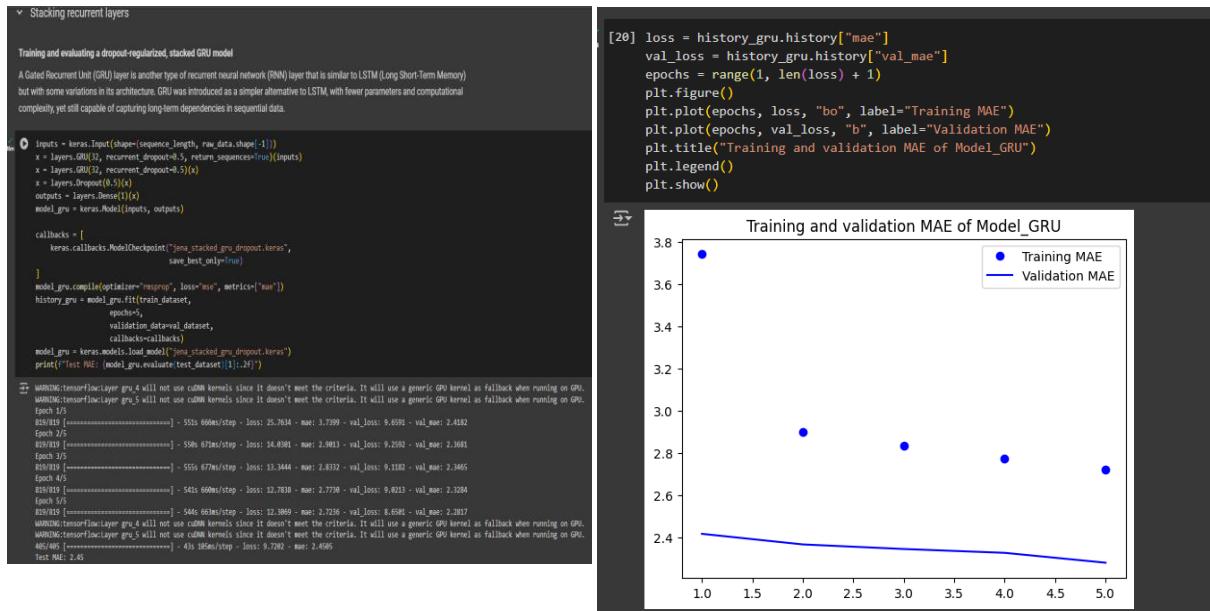
The model was trained for only 5 epochs due to computational limitations. Despite the limited training duration, the model achieved a test accuracy of MAE 2.56.



## Exploring GRU Network with Recurrent Dropout

Following the initial investigation with a basic model, we explored the utilization of a Gated Recurrent Unit (GRU) network for time series forecasting. This model incorporated two GRU layers with recurrent dropout applied at a rate of 0.5. Additionally, a dropout layer with a rate of 0.5 was included.

Due to computational limitations associated with recurrent dropout requiring a stronger GPU, the model was trained for only 5 epochs. Despite the limited training duration, the model achieved a promising validation Mean Absolute Error (MAE) of 2.45, with the potential for further improvement with additional training epochs.



## Week 10 and 11

### B. Analyzing Textual Time Travel Data with RNN Models

This report details the exploration of Recurrent Neural Networks (RNNs) for analyzing a unique dataset derived from "time travel data" stored in a .txt format.

#### Data Preprocessing:

- Tokenization:** The text data was split into lines, words, and then further divided into individual characters.
- Vocabulary Building:** Unique words were identified, resulting in a vocabulary size of 45,800.
- Corpus Creation:** The processed text was converted into a corpus suitable for RNN training.

```
[9]: from utils import *
import math
import torch
from torch import nn
from torch.nn import functional as F

[10]: torch.cuda.get_device_name(0)
[10]: 'NVIDIA GeForce RTX 3050 Ti Laptop GPU'

[11]: def tokenize(lines, token='word'):
        """Split text lines into word or character tokens."""
        if token == 'word':
            return [line.split() for line in lines]
        elif token == 'char':
            return [list(line) for line in lines]
        else:
            print('ERROR: unknown token type: ' + token)

[12]: def download(name, cache_dir=os.path.join('data')):
        """Download a file inserted into DATA_HUB, return the local filename."""
        assert name in DATA_HUB, f"Name {name} does not exist in {DATA_HUB}."
        url, sha1_hash = DATA_HUB[name]
        mkdir_if_not_exist(cache_dir)
        fname = os.path.join(cache_dir, url.split('/')[-1])
        print(f"File location: {fname}")
        if os.path.exists(fname):
            sha1 = hashlib.sha1()
            with open(fname, 'rb') as f:
                while True:
                    data = f.read(1048576)
                    if not data:
                        break
                    sha1.update(data)
            if sha1.hexdigest() == sha1_hash:
                return fname # Hit cache
        print(f'Downloading {fname} from {url}...')
        r = requests.get(url, stream=True, verify=True)
        with open(fname, 'wb') as f:
            f.write(r.content)

# We first tokenize text lines into words
tokens = tokenize(lines, token='word')

# Then, we construct vocabulary with the tokens
vocab = Vocab(tokens)

# Print the first few frequent tokens in the vocab with their indices.
print(f'vocab (size:{len(vocab)}): {list(vocab.token_to_idx.items())[:10]}')

# We can convert each text line into a list of numerical indices
for i in [0, 10]:
    print('words:', tokens[i])
    print('indices:', vocab[tokens[i]])

vocab (size:4580): [('unk', 0), ('the', 1), ('i', 2), ('and', 3), ('of', 4),
words: ['the', 'time', 'machine', 'by', 'h', 'g', 'wells']
indices: [1, 19, 50, 40, 3130, 3058, 438]
words: ['twinkled', 'and', 'his', 'usually', 'pale', 'face', 'was', 'flushed'
indices: [4399, 3, 25, 1398, 387, 113, 7, 1676, 3, 1053, 1]

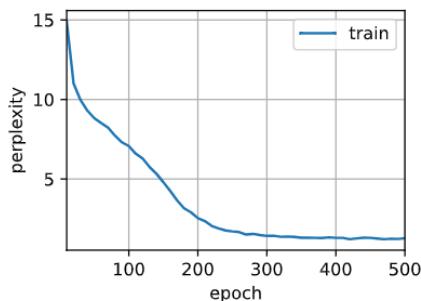
corpus, vocab = load_corpus_time_machine()
len(corpus), len(vocab)
File location: data\timemachine.txt
(170580, 28)
```

## Model Training and Evaluation:

Several RNN architectures were evaluated on the time travel data corpus with a batch size of 32 and a sequence length of 35 characters (Num\_steps).which determines the length of each sequence processed by the RNN during training. The models were trained for 500 epochs, and the following results were obtained:

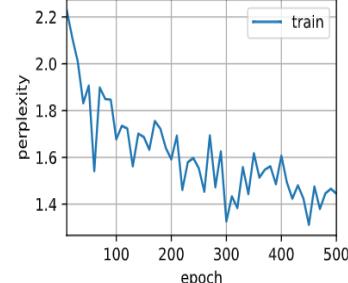
- **Baseline RNN:**
  - Perplexity: 1.3 (Lower perplexity indicates better performance)
  - Token Processing Speed: 39,866.4 tokens/sec
  - **Random Iteration:** Training with random iteration resulted in a slight decrease in performance (perplexity: 1.4).
- **Concise RNN:**
  - This model achieved a significant improvement in processing speed:
    - Perplexity: 1.3 (similar to Baseline RNN)
    - Token Processing Speed: 497,418.5 tokens/sec (over 12 times faster)
  - **Random Iteration:** Here, a surprising trend emerged. Training with random iteration led to a slight improvement in perplexity (1.2) while maintaining a high processing speed (557,043.9 tokens/sec).

```
num_epochs, lr = 500, 1
train_rnn(rnn_scratch, train_iter, vocab, lr, num_epochs, try_gpu())
perplexity 1.3, 98664.9 tokens/sec on cuda:0
time travelleryou can show black is white by argument said filby
traveller hitheress boted one mithtmere washad he this ingl
```



Vanilla RNN

```
train_rnn(rnn_scratch, train_iter, vocab, lr, num_epochs, try_gpu(),
use_random_iter=True)
perplexity 1.4, 93466.6 tokens/sec on cuda:0
time travellerit s against reason said the time travellerit s ag
travellerit s against reason said the time travellerit s ag
```



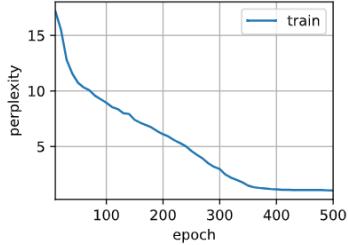
Concise RNN

## Advanced RNN Architectures:

- **GRU (Gated Recurrent Unit):**
  - **Scratch Model:** This model achieved the best perplexity score (1.1) but with a slower processing speed (31,268.5 tokens/sec).
  - **Concise Model:** This model offered a balance between performance and speed, achieving the lowest perplexity (1.0) among all models while maintaining a reasonable processing speed (320,064.1 tokens/sec).

```
vocab_size, num_hiddens, device = len(vocab), 256, try_gpu()
num_epochs, lr = 500, 1
gru_scratch = RNNModelScratch(len(vocab), num_hiddens, device, init_gru_params,
                               init_gru_state, gru)
train_rnn(gru_scratch, train_iter, vocab, lr, num_epochs, device)
```

perplexity 1.1, 31268.5 tokens/sec on cuda:0  
time traveller for so it will be convenient to speak of himwas e  
traveller with a slight accession of cheerfulness really thi

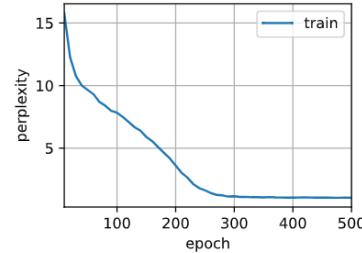


GRU\_Scratch

### GRU Concise

```
]: num_inputs = vocab_size
gru_layer = nn.GRU(num_inputs, num_hiddens)
gru_concise = RNNModel(gru_layer, len(vocab))
gru_concise = gru_concise.to(device)
train_rnn(gru_concise, train_iter, vocab, lr, num_epochs, device)
```

perplexity 1.0, 320064.1 tokens/sec on cuda:0  
time traveller you can show black is white by argument said filby  
traveller you can show black is white by argument said filby

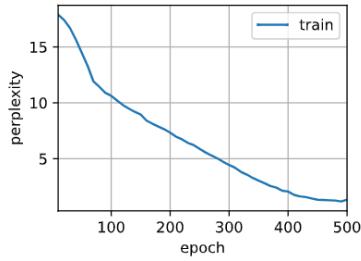


Concise GRU

- **LSTM (Long Short-Term Memory):**
  - **Standard Model:** Similar to the Baseline RNN, the standard LSTM achieved an average perplexity (1.3) but was slower (27,329.0 tokens/sec).
  - **Concise Model:** This model achieved the same best perplexity (1.0) as the concise GRU but with a slight processing speed advantage (248,827.8 tokens/sec).

```
]: vocab_size, num_hiddens, device = len(vocab), 256, try_gpu()
num_epochs, lr = 500, 1
lstm_scratch = RNNModelScratch(len(vocab), num_hiddens, device, init_lstm_params,
                               init_lstm_state, lstm)
train_rnn(lstm_scratch, train_iter, vocab, lr, num_epochs, device)
```

perplexity 1.3, 27329.0 tokens/sec on cuda:0  
time traveller cullenthe simell be cofisars of sace as i bone o  
traveller cullesse which i will explais of space but you ca

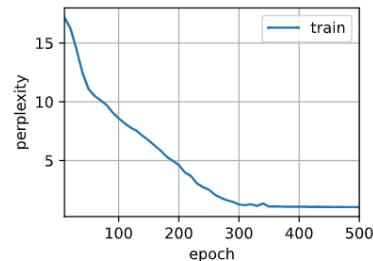


LSTM\_Scratch

### LSTM Concise

```
]: num_inputs = vocab_size
lstm_layer = nn.LSTM(num_inputs, num_hiddens)
lstm_concise = RNNModel(lstm_layer, len(vocab))
lstm_concise = lstm_concise.to(device)
train_rnn(lstm_concise, train_iter, vocab, lr, num_epochs, device)
```

perplexity 1.0, 248827.8 tokens/sec on cuda:0  
time traveller for so it will be convenient to speak of himwas e  
traveller you can show black is white by argument said filby



## Conclusion

- This project demonstrates the effectiveness of using RNN, GRU, and LSTM architectures for sequence modeling tasks on text data. The concise implementations of these models not only improve training speed significantly but also maintain or enhance the model's performance as measured by perplexity.
- **GRU and LSTM models** generally perform better than vanilla RNNs in terms of perplexity, making them more suitable for complex sequence modeling tasks.
- **Concise models** (pre-built) offer substantial computational benefits, achieving higher throughput in tokens per second while maintaining competitive perplexity scores.