

# Chapter 2 - Compression Techniques

*"I have made this longer than usual because I have not had time to make it shorter."*

*Blaise Pascal*

In the last chapter, we discussed a few ideas to improve the deep learning efficiency. Now, we will elaborate on one of those ideas, the compression techniques. Compression techniques aim to reduce the model footprint (size, latency, memory etc.). We can reduce the model footprint by reducing the number of trainable parameters. However, this approach has two drawbacks. First, it is hard to determine the parameters or layers that can be removed without significantly impacting the performance. It requires many trials and evaluations to reach a smaller model, if it is at all possible. Second, such an approach doesn't generalize well because the model designs are subjective to the specific problem.

In this chapter, we introduce Quantization, a model compression technique that addresses both these issues. We'll start with a gentle introduction to the idea of compression. Details of quantization and its applications in deep learning follow right after. The quantization section delves into the implementation details using code samples. We finish with a hands-on project that will walk you through the process of applying quantization in practical situations using popular frameworks like Tensorflow and Tensorflow Lite.

## An Overview of Compression

One of the simplest approaches towards efficiency is compression to reduce data size. For the longest time in the history of computing, scientists have worked tirelessly towards storing and transmitting information in as few bits as possible. Depending on the use case, we might be interested in compressing in a lossless or lossy manner. We can fit 10 apples in a smaller box with a better arrangement. This is lossless compression. Another approach is to chop them into cubes and discard the odd parts. An even smaller box can fit those 10 apples this way. We can call this lossy compression because we lost the odd parts. The choice of the technique depends on several factors like customer preference, consumption delay, or resource availability (extra hands needed for chopping). Personally, I like *full* apples.

Let's move on from apples to the digital domain. A popular example of lossless data compression algorithm is [Huffman Coding](#), where we assign unique strings of bits (codes) to the symbols based on their frequency in the data. More frequent symbols are assigned smaller codes, and less frequent symbols are assigned longer codes. This is achieved with a simple Huffman Tree (figure 2-1 bottom). Each leaf node in the tree is a symbol, and the path to that symbol is the bit-string assigned to it. This allows us to encode the given data in as few bits as possible, since the most frequent symbols will take the least number of bits to represent. In aggregate, this would be better than encoding each symbol with the same number of bits. The lookup table (figure 2-1 middle) that contains the symbol-code mapping is transmitted along with the encoded data.

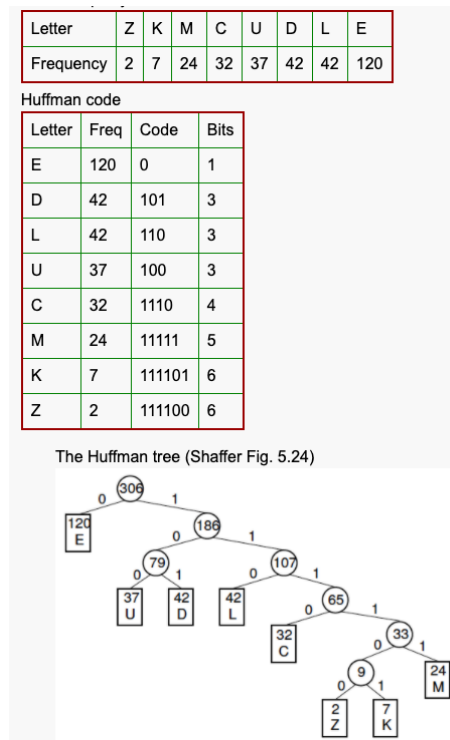


Figure 2-1: Huffman Encoding & Huffman Tree. [Source](#)

When decoding the encoded data, we look up the code from the lookup table to retrieve the symbols back. Since the codes are unique for each symbol (in fact, they are prefix codes: no code is a prefix of some other code, which eliminates ambiguity when decoding), we can easily construct the original sequence of symbols from the encoded sequence and the lookup table. Refer the wikipedia article on [arithmetic coding](#) to learn about lossless coding schemes.

The lossy compression algorithms are used in situations (people who like diced apples) where we don't expect to recover the exact representation of the original data. It is okay to recover an approximation, however we do expect a better compression ratio than lossless compression, since we are losing some information as a trade off. It is especially applicable for multimedia (audio, video, images) data,, where it is likely that either humans who will consume the information will not notice the loss of some information, or do not necessarily care about the loss in quality.



Figure 2-2: On the left is a high quality image of a cat. The cat on the right is a lower quality compressed image. [Source](#)

Both the cat images in figure 2-2 might serve their purpose equally well, but the compressed image is an order of magnitude smaller. Discrete Cosine Transform (DCT), is a popular algorithm which is used in the JPEG format for image compression, and the MP3 format for audio. DCT breaks down the given input data into independent components of which the ones that don't contribute much to the original input can be discarded, based on the tolerance for loss in quality. The JPEG and MP3 formats are able to achieve a 10-11x compression without any perceptible loss in quality. However, further compression might lead to degradation in quality.

In our case, we are concerned about compressing the deep learning models. What do we really mean by compressing though? As mentioned in chapter 1, we can break down the metrics we care about into two categories: *footprint metrics* such as model size, prediction latency, RAM consumption and the *quality metrics*, such as accuracy, F1, precision and recall as shown in table 2-1.

Footprint Metrics	Quality Metrics
<ul style="list-style-type: none"><li>• Model Size</li><li>• Inference Latency on Target Device</li><li>• Training Time for Convergence</li><li>• Peak RAM Consumption</li></ul>	<ul style="list-style-type: none"><li>• Accuracy</li><li>• Precision</li><li>• Recall</li><li>• F1</li><li>• AUC</li></ul>

Table 2-1: A few examples of footprint and quality metrics.

The footprint and the quality metrics are typically at odds with each other. As stated earlier for JPEG and MP3 encoding, compression beyond a limit hurts quality metrics. Conversely, a higher quality implies a worse footprint. In the case of deep learning models, the model quality is often correlated with the number of layers, and the number of parameters (assuming that the models are well-tuned). If we naively reduce the footprint, we can reduce the number of layers and number of parameters, but this could hurt the quality.

Compression techniques are used to achieve an efficient representation of one or more layers in a neural network with a possible quality trade off. The efficiency goals could be the optimization of the model with respect to one or more of the footprint metrics such as the model size, inference latency, or training time required for convergence with a little quality compromise. Hence, it is important to ensure that we evaluate these techniques using the metrics relevant to our use case. In some cases, these techniques can help reduce complexity and improve generalization.

Let us consider an arbitrary neural network layer. We can abstract it using a function  $f$  with an input  $x$  and parameters  $\theta$  such that  $y = f(x; \theta)$ . In the case of a fully-connected layer,  $\theta$  is a 2-D matrix. Further, assume that we can train another network  $y' = f'(x; \theta')$  with far fewer parameters ( $\theta' \ll \theta$ ) such that the outputs are approximately the same  $y' \approx y$ . Such a model  $f'$  is useful if we want to deploy a model in a space constrained environment like a mobile device.

To summarize, compression techniques help to achieve an efficient representation of a layer or a collection of layers, such that it meets the desired tradeoff goals. In the next section we introduce Quantization, a popular compression technique which is also used in various fields of computer science in addition to deep learning.

## Quantization

Before we jump to working with a deep learning model, we have a task for you. You have been handed the charge of the Mars Rover! The rover is transmitting images back to earth. However, transmission costs make it infeasible to send the original image.

Can we compress the transmission, and decompress it on arrival? If so, what would be the ideal tradeoff on how much compression we want v/s how much quality loss can we tolerate?

Let us slowly build up to that by exploring how quantization can help us.

### A Generic View of Quantization

Quantization is a common compression technique that has been used across different parts of Computer Science especially in signal processing. It is a process of converting high precision continuous values to low precision discrete values. Take a look at figure 2-3. It shows a sine wave and an overlapped quantized sine wave. The sine wave is continuous, a high precision representation. The quantized sine wave is a low precision representation which takes integer values in the range  $[0, 5]$ . As a result, the quantized wave requires low transmission bandwidth.

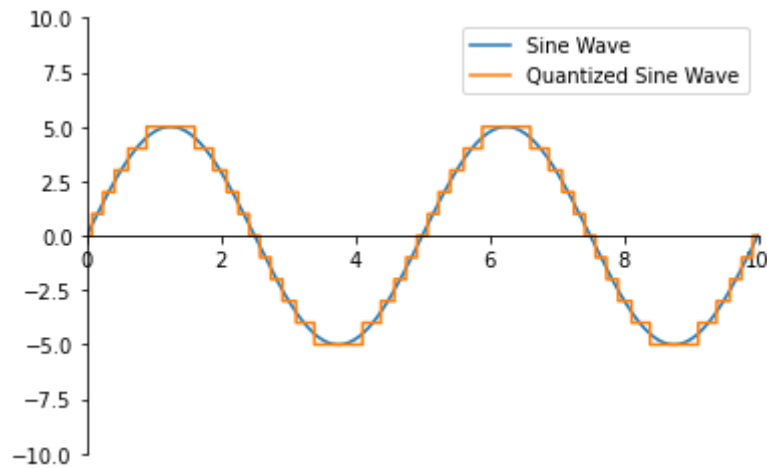


Figure 2-3: Quantization of sine waves.

Let's dig deeper into its mechanics using an example. Let's assume we have a variable  $x$  which takes a 32-bit floating point value in the range  $[-10.0, 10.0]$ . We need to transmit a collection (vector) of these variables over an expensive communication channel. Can we use quantization to reduce transmission size and thus save some costs?

What if it did not matter to us if  $x$  was stored/transmitted with some error ( $-5.023$  v/s  $-5.0$ )? If we can tolerate some loss of precision, can we use  $b$ -bits and save some space? Let us work on a scheme for going from this higher-precision domain (32-bits) to a quantized domain ( $b$ -bit values). This process is nothing but (cue drum roll!) ...Quantization!

Before we get our hands dirty, let us first make two reasonable assumptions:

1. We know that the value of  $x$  will lie between  $-10.0$  ( $x_{\min}$ ) and  $10.0$  ( $x_{\max}$ ). Let us assume that this will always hold true. If not, the value of  $x$  will be clamped to lie in this range.
2. Let us assume that the values of  $x$  will be uniformly distributed in this range. This means that all values of  $x$  are equally likely to lie in any part of the range from  $[x_{\min}, x_{\max}]$ , and there are no clusters of values in any part.

Now that we have the assumptions out of the way, instead of working with a 32-bit for storing  $x$ , let us assume we have a  $b$ -bit unsigned integer for storing  $x$ . A  $b$ -bit unsigned integer will have  $2^b$  possible distinct values, ranging from 0 to  $2^b - 1$ .

To go from a 32-bit floating point value to a  $b$ -bit integer, and back again, we need a mapping from one side to the other. It is easy to learn a mapping from 32-bit to  $b$ -bit values.

We would also want to keep a weaker relative ordering of the values between the two domains. That is, if  $a < b$  in the higher-precision domain,  $a \leq b$  in the quantized domain. We have the  $\leq$  in the low precision domain, because we are losing precision when going to a  $b$ -bit integer and as a result values which were close in the high precision domain might end up being mapped to the same value. For example,  $-10.0$  and  $-9.999$  might both be mapped to 0 in the quantized domain.

Keeping all that in mind, it is easy to see that floating-point  $x_{\min}$  should map to 0, and  $x_{\max}$  should map to  $2^b - 1$ . How do we map the rest of the floating point values in between  $x_{\min}$  and  $x_{\max}$  to integer values?

### Exercise: Mapping from a high precision to a low precision domain.

Visually inspecting figure 2-4, can you work out the formula for mapping a given floating-point value ( $x$ ) to a quantized value ( $x_q$ ). Assume that you are given values of  $x_{\min}$ ,  $x_{\max}$ , and  $b$ ?

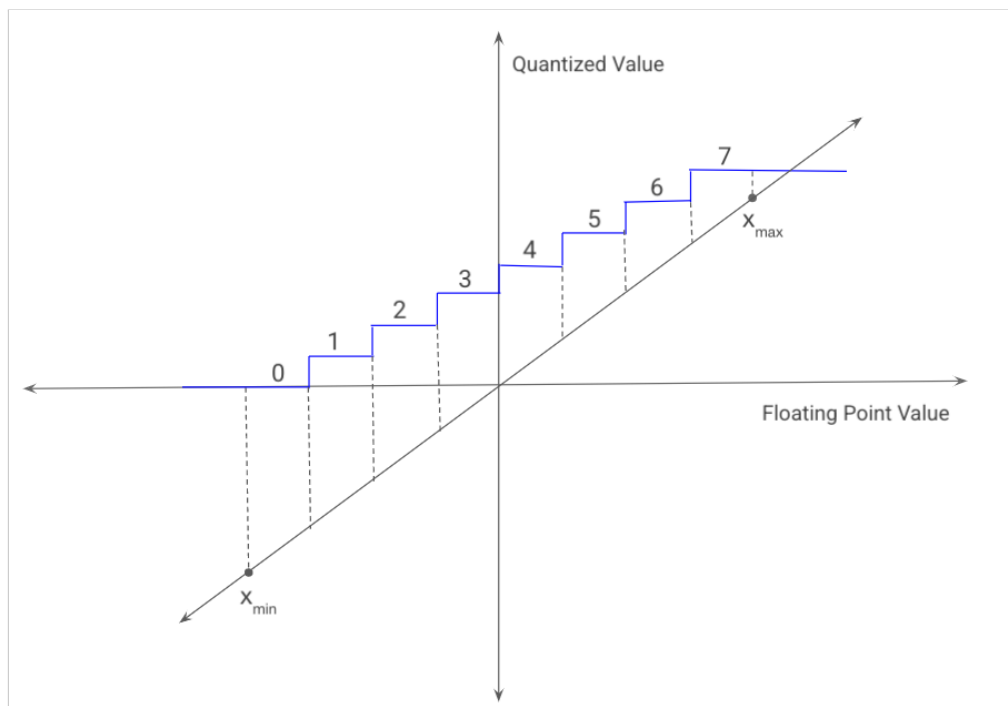


Figure 2-4: Quantizing floating-point continuous values to discrete unsigned values. The continuous values range from  $x_{\min}$  to  $x_{\max}$ , and are mapped to continuous values in  $[0, 2^b - 1]$  (in the above figure,  $b = 3$ , hence the quantized values are in the range  $[0, 7]$ ). For the purpose of quantization, the continuous values are also clamped to be in the range  $[x_{\min}, x_{\max}]$ .

#### Solution:

Note that we have to map all the values from  $[x_{\min}, x_{\max}]$  to  $2^b$  possible values (let's call them bins). Figure 2-5 shows a visual representation of the mapping. The values of  $x$  are uniformly spread out (as marked using ticks in figure 2-5). Hence, every bin in the quantized domain should have an equal sized range in the higher-precision domain mapping towards it.

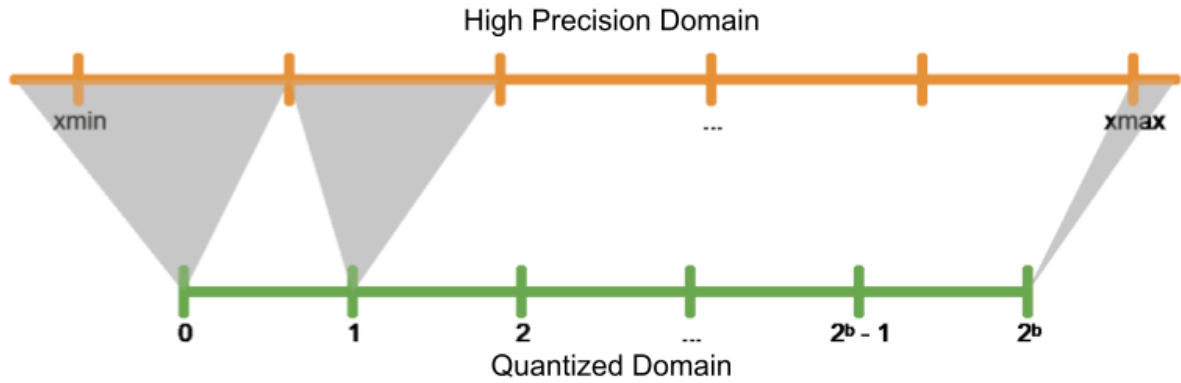


Figure 2-5: Mapping from  $[x_{\min}, x_{\max}]$  to the quantized domain.

Since there are  $2^b$  possible bins, to be divided equally amongst the range  $[x_{\min}, x_{\max}]$ , each bin covers a range of  $s = \frac{x_{\max} - x_{\min}}{2^b - 0} = \frac{x_{\max} - x_{\min}}{2^b}$  (where,  $s$  is referred to as scale).

Hence,  $[x_{\min}, x_{\min} + s)$  will map to the bin 0,  $[x_{\min} + s, x_{\min} + 2s)$  will map to bin 1,  $[x_{\min} + 2s, x_{\min} + 3s)$  will map to bin 2, and so on.

Thus, to find which bin the given  $x$  will go to, we simply do the following:

$$x_q = \left\lfloor \frac{x - x_{\min}}{s} \right\rfloor.$$

We need the floor function ( $\lfloor \cdot \rfloor$ ) so that the floating point value is converted to an integer value.

Now, if we plug in  $x = x_{\min}$ , and  $x_q$  would be 0.

Although, when we plug in  $x = x_{\max}$ ,  $x_q$  would be  $2^b$ , which is not okay since the maximum value in the quantized domain is  $2^b - 1$ .

So we would clamp it back to fit in the range  $[0, 2^b - 1]$ . To summarize, our formula becomes:

$$x_q = \min(x_q, 2^b - 1),$$

$$\text{where } s = \frac{x_{\max} - x_{\min}}{2^b}.$$

In this exercise, we worked out the formulas to map a high precision domain to a low precision domain. The following exercise will apply them to quantize an arbitrary data sequence.

### Exercise: Data Quantization

Let's put our learnings from the previous exercise into practice. We will code a method `quantize` that quantizes a vector  $x$ , given  $x_{\min}$ ,  $x_{\max}$ , and  $b$ . It should return the quantized values for a given  $x$ .

### Solution:

We use [NumPy](#) for this solution. It supports vector operations which operate on a vector (or a batch) of  $x$  variables (vectorized execution) instead of one variable at a time. Although it is possible to work without it, you would have to introduce a for-loop either within the function, or outside it. This is crucial for deep learning applications which frequently operate on batches of data. Using vectorized operations also speeds up the execution (and this book is about efficiency, after all!). We highly recommend learning and becoming familiar with numpy.

```
# numpy is one of the most useful libraries for ML.
import numpy as np

def get_scale(x_min, x_max, b):
    # Compute scale as discussed.
    return (x_max - x_min) * 1.0 / (2**b)

"""Quantizing the given vector x."""
def quantize(x, x_min, x_max, b):
    # Clamp x to lie in [x_min, x_max].
    x = np.minimum(x, x_max)
    x = np.maximum(x, x_min)

    # Compute scale as discussed.
    scale = get_scale(x_min, x_max, b)
    x_q = np.floor((x - x_min) / scale)

    # Clamping the quantized value to be less than (2^b - 1).
    x_q = np.minimum(x_q, 2**b - 1)

    # Return x_q as an unsigned integer.
    # uint8 is the smallest data type supported by numpy.
    return x_q.astype(np.uint8)
```

The code is self-explanatory. We receive a vector of  $x$  values, and then we clamp it in the  $[x_{\min}, x_{\max}]$  range. Then we compute the scale as discussed previously, which is used to compute  $x_q$ .  $x_q$  is also clamped to ensure that it does not exceed  $2^b - 1$ .

The final returned vector's type is converted to the uint8 data type. Note that  $b$  might be less than 8, in which case uint8 leads to unnecessary space wastage. If that is indeed the case, you might have to design your own mechanism to pack in multiple quantized values in one of the supported data types (using bit-shifting).

For example, if you pick  $b = 4$ , you might want to pack two quantized values in a uint8 manually, and unpack them when decoding the data.

Now let's run the code for the range  $[-10, 10]$ , incrementing by 2.5 each time and find the quantized values for  $b = 3$ . First, let's create our  $x$ .



```
# Construct the array that we wish to quantize.
# We slightly exceed 10.0 to include 10.0 in our range.
x = np.arange(-10.0, 10.0 + 1e-6, 2.5)
print(x)
```

We do this using NumPy's `arange` method, which allows us to generate a range of floating point values, with the starting and endpoint defined, along with a step value. This returns the following result.

```
[-10.   -7.5   -5.    -2.5    0.     2.5    5.     7.5   10. ]
```

Now let's quantize `x`.

```
# Quantize the entire array in one go.
x_q = quantize(x, -10.0, 10.0, 3)
print(x_q)
```

This returns the following result.

```
[0 1 2 3 4 5 6 7 7]
```

Table 2-2 shows the element wise comparison of `x` and `xq`.

<b>x</b>	-10.0	-7.5	-5.0	-2.5	0	2.5	5.0	7.5	10.0
<b>x<sub>q</sub></b>	0	1	2	3	4	5	6	7	7

Table 2-2: Vector `x` and the quantized vector `xq`

As we had calculated earlier, the value `r` was 2.5 for our case. Hence the range `[-10.0, -7.5)` maps to 0, `[-7.5, -5.0)` maps to 1, and so on. `[-7.5, 10.0)` maps to 7. Since 10.0 is just outside the range, but the maximum possible value of `xq` is 7, it is clamped to that value.

In the last two exercises, we worked out the logic to quantize a high precision vector to low precision to save storage space and the transmission bandwidth. Let's say a receiver received this data. How would it decode it to get the original value? The next exercise details a dequantization algorithm to retrieve the original (with an acceptable tolerance) value.

## Exercise: Data Dequantization

*"But you wouldn't clap yet. Because making something disappear isn't enough; you have to bring it back. That's why every magic trick has a third act, the hardest part, the part we call 'The Prestige'."*

The Prestige (2006)

Quantization is just half a piece of the puzzle. We need to be able to bring back the original value (with an acceptable tolerance). Given our work so far, let's try dequantizing the encoded value?

### Solution:

For dequantization, we compute the scale the same way as earlier.  $x_q$  is our bin number, which along with the scale can tell us how far from the  $x_{min}$  we are.

Remember that  $[x_{min}, x_{min} + s)$  will map to the bin 0,  $[x_{min} + s, x_{min} + 2s)$  will map to bin 1,  $[x_{min} + 2s, x_{min} + 3s)$  will map to bin 2, and so on. The values in bin 0 in the quantized domain will map to  $x_{min}$  in the high precision domain. Similarly, bin 1 values in the quantized domain will map to  $x_{min} + s$ . And so on. Hence the dequantized value of  $x$  would be

$$x = x_{min} + sx_q.$$

```
def dequantize(x_q, x_min, x_max, b):  
    # Compute the value of scale the same way.  
    s = get_scale(x_min, x_max, b)  
    x = x_min + (s * x_q)  
    return x
```

Let's dequantize<sup>1</sup> the quantized vector  $x_q$ .

```
dequantize(x_q, -10.0, 10.0, 3)
```

We receive this dequantized array upon running the code. Note that the last element was supposed to be 10.0, and the error is 2.5.

```
array([-10. , -7.5, -5. , -2.5,  0. ,  2.5,  5. ,  7.5,  7.5])
```

So far so good! We have learnt to dequantize a quantized vector. Let's keep going and apply these ideas on a regular *jpeg* image in the next exercise.

### Exercise: Quantize an image, then dequantize it.

Let's go back to the original task of optimizing the Mars rover's communications! Can we quantize the transmission, and dequantize it on arrival? If so, what would be the ideal number of bits we can use so the quality does not degrade too much? The original (pre-quantization) image is shown in figure 2-6.

Get the image using this command:

```
!wget https://github.com/reddragon/book-codelabs/raw/main/pia23378-16.jpeg
```

### Solution:

---

<sup>1</sup> To dequantize, we map the quantized value to the smallest value in the dequantized range of the bin. For example, quantized 0 will map to  $x_{min}$ . There could be other schemes mapping schemes as well such as mapping it to the middle of the range, or to the maximum value in the range. However, our approach is what is used in deep learning conventionally.

First, we will interpret the image in the form of a 2D matrix having values in [0.0, 1.0].

```
%matplotlib inline

import matplotlib.pyplot as plt
import matplotlib.image as mpimg

img = (mpimg.imread('pia23378-16.jpg') / 255.0)
plt.imshow(img)
```



Figure 2-6: Mars Curiosity Rover. Dimensions: 586x1041. [Source](#)

Now, let's simulate the process of transmission by quantization, followed by dequantization. We will reuse the *quantize()* and *dequantize()* methods from the previous exercises. However, in this case, the image data values lie in range [0.0, 1.0]. Take a look at the *simulate\_transmission()* method below which uses this range to call *quantize()* and *dequantize()* methods.

```
def simulate_transmission(img, b):
    transmitted_image = quantize(img, 0.0, 1.0, b)
    decoded_image = dequantize(transmitted_image, 0.0, 1.0, b)
    plt.axis('off')
    plt.imshow(decoded_image)
```

Figure 2-7 shows quantized images for different values of *b*. The image in the top left corner (*b*=2) uses 2-bits to represent each value. The quality degradation is apparent. The image in the bottom right corner (*b*=8) is the highest quality image. It's hard to notice much difference for the *b* values 6, 7, and 8. The ideal spot seems to be at *b*=5. Since the original image was 8-bit encoded, with a 5-bit encoding we save 38.5% space while still getting a reasonable quality image.

b=2	b=4
-----	-----

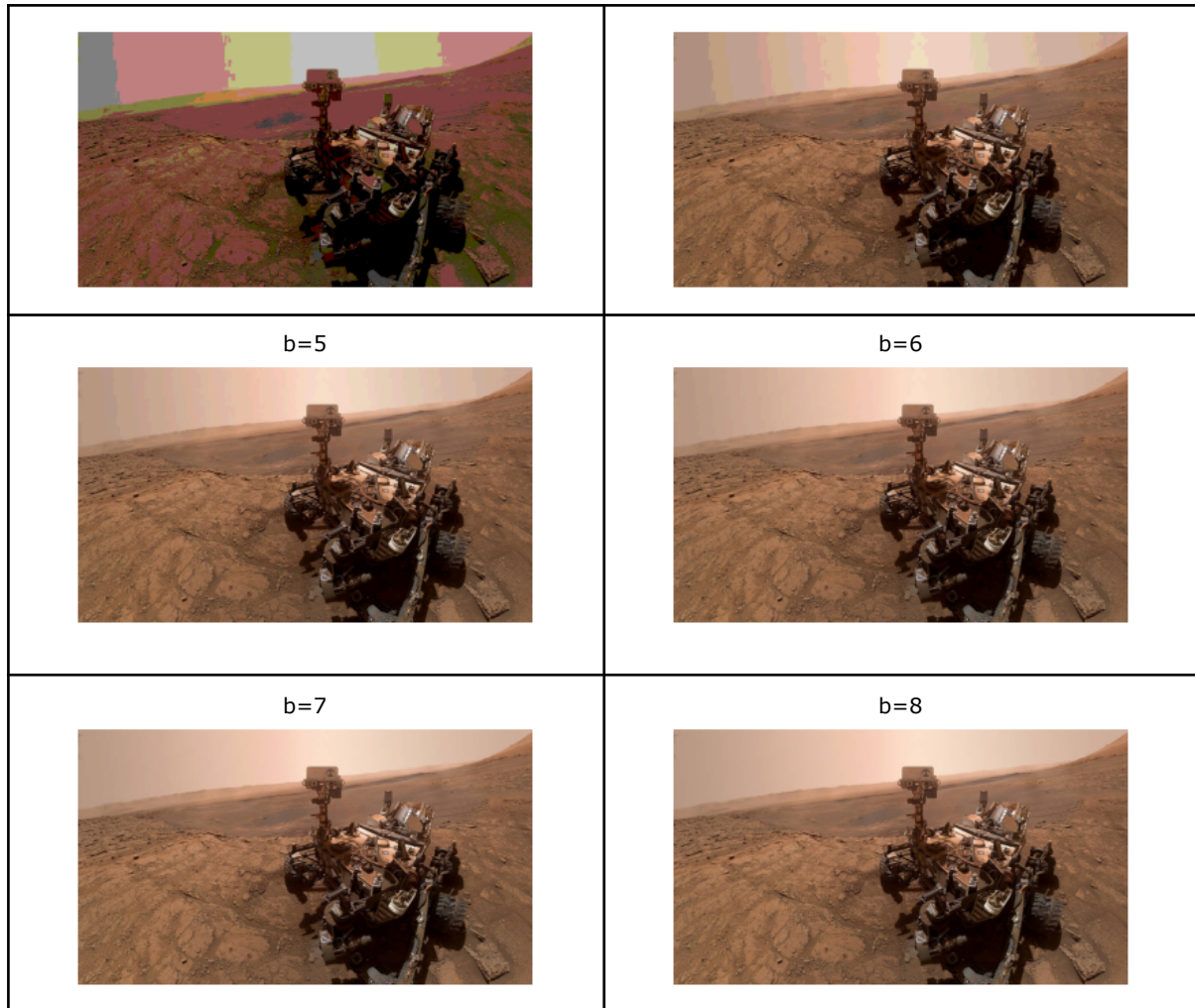


Figure 2-7: Images with various degrees of quantization.

A graph of quantized representation bit size ( $b$ ) and the resulting image sizes (in bits) is shown in figure 2-8. It demonstrates that the sizes of the resulting image drop linearly with the reduction in the number of quantization bits. Quantization is a useful technique in the situation where the storage space or the transmission bandwidth is expensive like deep learning models on mobile devices. Mobile devices are resource constrained. Hence, quantization can help to deploy models which would otherwise be too big to execute on such devices. We will tackle this exact problem in the following section.

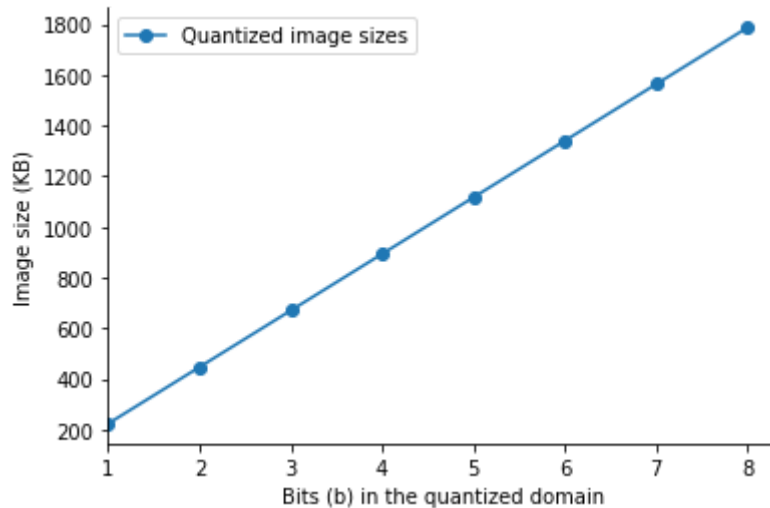


Figure 2-8: Image sizes with various degrees of quantization

## Quantization in Deep Learning Models

It was exciting to compress the transmission of the Mars Rover! Now, how can we use quantization for deep learning models?

Most of the model size of a Deep Learning model comes from the weights in its layers. Similarly, most of the latency comes from computing the activations. Typically, the weights and activations are 32-bit floating-point values. One of the ideas for reducing model footprint is to reduce the precision for the weights and activations, by quantizing them into a lower-precision data type (often 8-bit integers), similar to what we just learnt.

There are two kinds of gains that we can get from quantization in Deep Learning models:

- (a) lower model size, and
- (b) lower inference latency.

We already have the necessary tools for achieving (a), the lower model size. Let us see how we can apply what we learnt for quantizing deep learning models.

## Looking Under the Hood

As we know, one of the basic neural network operation is as follows:

$$f(X; W, b) = \sigma(XW + b)$$

Here,  $X$ ,  $W$  and  $b$  are tensors (mathematical term for an  $n$ -dimensional matrix) to denote inputs, weights and the bias respectively. To simplify this discussion, let's assume the shape (an array describing the size of each dimension) of  $X$  as  $[\text{batch size}, D_1]$ , that of  $W$  as  $[D_1, D_2]$  and  $b$  is the bias vector with shape  $[D_2]$ .

Hence, the shape of the result of the operation  $(XW + b)$  is  $[\text{batch size}, D_2]$ .  $\sigma$  is a nonlinear function that is applied element-wise to the result of  $(XW + b)$ . Some examples of the nonlinear functions are ReLU ( $\text{ReLU}(x) = x$  if  $x > 0$ , else 0), tanh, sigmoid, etc.

A neural network model learns  $W$  and  $b$  tensors which are stored with the model. Hence, they contribute significantly to its size. Of the two tensors,  $W$  naturally dominates the size owing to its higher dimensionality. Needless to say, an efficient representation of  $W$  is necessary to reduce its size.

In terms of latency, the most expensive operation is  $XW$  which is a matrix<sup>2</sup> multiplication operation. The next operation  $(XW + b)$  is a vector addition and  $\sigma$  is an element-wise operation. Both of these operations are cheaper to compute than matrix-multiplication. To optimize the computation latency, we should compute  $XW$  in an efficient manner.

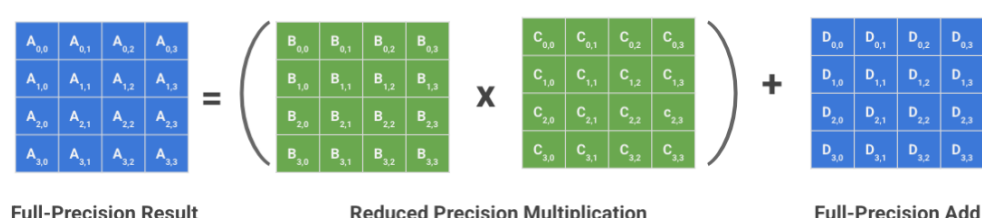


Figure 2-9: An optimized Multiply-Accumulate Operation. BC is the most expensive operation latency-wise. In typical neural-networks:  $C$  is the weight matrix.  $D$  is often a one-dimension vector, hence the addition is cheap both from the latency point of view and size wise (since  $C$  dominates the size).

In fact, the general formulation of  $Y = XW + b$ , is the [Multiply-Accumulate operation](#) (MAC). Figure 2-9 describes the MAC operation for  $A = BC + D$ .  $B$ ,  $C$ , and  $D$  are all matrices.

The number of MACs in a model is another metric to measure its complexity. And since they are such a fundamental block of models, they have been optimized both in hardware and software. Let's take a look at how we can optimize a slightly easier version of this operation (where  $D$  is a vector instead of a matrix) using quantization.

## Weight Quantization

We just learnt that the weights in deep learning models are stored in an  $N$ -dimensional matrix (tensor), and the weight matrix  $W$  is most expensive in terms of storage. Can we efficiently represent this weight matrix  $W$  to reduce the model size? We already have worked with 1-D vectors and 2-D vectors in previous exercises. We can extend that learning to deep learning matrices as well.

Here is a simple scheme for quantizing weights to reduce a model's size:

<sup>2</sup>The terms tensor and matrix are used interchangeably in the text. Tensors are  $N$ -dimensional matrices. The shape of a tensor denotes the size of each of its dimensions, and the rank of the tensor denotes the number of dimensions of that tensor.

1. Given a 32-bit floating-point weight matrix, we can map the minimum weight value  $x_{\min}$  to 0, and the maximum weight value  $x_{\max}$  to  $2^b - 1$  ( $b$  is the number of bits of precision and  $b < 32$ ). Notice how this is similar to computing the  $x_{\min}$  and  $x_{\max}$  of an arbitrary matrix.
2. Then we can map all the values in the weight matrix to an integer value in range  $[0, 2^b - 1]$ . This is quantization of the weight matrix by mapping a floating point value to a fixed-point value where the latter requires a lesser number of bits.
3. This process can also be applied to signed  $b$ -bit fixed-point integers, where the output values will be in the range  $[-2^{\frac{b}{2}} - 1, 2^{\frac{b}{2}} - 1]$ . One of the reasonable values of  $b$  is 8, since this would lead to a  $32 / 8 = 4x$  reduction in space. This fits in well since there is near-universal support for unsigned and signed 8-bit integer data types.
4. The quantized weights are persisted with the model.

With this scheme, we can shrink the model sizes with an acceptable loss of precision. A smaller model size can be deployed in resource constrained environments like the mobile devices. Quantization has enabled a whole lot of models to run on mobile devices and IoTs which otherwise wouldn't be possible. We have solved one piece of the puzzle. Let's tackle the remaining piece, dequantization, next!

## Dequantization

It is not sufficient to just store the models. Our goal is to use them to make inferences (predict the output for a given input). Quantization transformed our weight matrices to a quantized integer format. However, the model layers use floating point computations. It is critical that our models are accurate in the inference stage.

Let's figure out how to dequantize these weight matrices. During inference, prior to the  $XW$  multiplication operation, we must dequantize  $W$ . It introduces a small computation latency. We can dequantize the weights using  $x_{\min}$ ,  $x_{\max}$ , and  $b$  just like we did in the data dequantization exercise. Note that the dequantization is a lossy process (figure 2-5) which recovers approximations of the original pre-quantization values. Once we have a dequantized approximation of the original matrix, we can proceed as usual.

However, we should expect an impact on the output quality based on how many bits were used for representing each value in the quantized domain.

## Exercise: Quantization simulation for a single fully connected layer.

Can you simulate the compression of a single fully connected layer using quantization? You can leverage the `np.random.uniform()` function (from the numpy package) to create dummy inputs ( $X$ ), weights ( $W$ ) and bias ( $b$ ) tensors. Using these three tensors, compute the layer output. Now, quantize the weights using  $b = 8$  and recompute the output using dequantized weights. How different are the two outputs?

### Solution:

We will start with the random number generator with a fixed seed to get consistent results across multiple runs. Next, we will create an input tensor of shape  $[10, 3]$ , where 10 is the



batch size, and 3 is the input dimension ( $D_1$  as stated earlier). The shape of the weights tensor is [3, 5], 3 is chosen to match  $D_1$ , and 5 is  $D_2$ . Bias is of shape [5]. The shapes are arbitrarily chosen for illustration purposes.

```
# Set the seed so that we get the same initialization.
np.random.seed(10007)

def get_random_matrix(shape):
    return np.random.uniform(low=-1.0, high=1.0, size=shape)

# Populate the inputs, weights and bias.
inputs = get_random_matrix([10,3])
weights = get_random_matrix([3,5])
bias = get_random_matrix([5])
```

Print the weights for comparison later on.

```
print(weights)
[[-0.08321415 -0.66657766  0.71264132 -0.39179407  0.05601718]
 [-0.85867389 -0.00864216 -0.15913464 -0.00676971  0.33190099]
 [-0.25760764 -0.82441528  0.57125625  0.74180458 -0.75251044]]
```

Let's calculate  $XW + b$ .

```
y = np.matmul(inputs, weights) + bias
print(y)
[[ 0.00511569 -0.89318885  0.51116489  0.57396818 -0.34144945]
 [ 1.34222938 -0.1270941  0.34179184 -0.75315659  0.39381581]
 [ 0.9916536  -1.15636837  0.96734553  0.53233659 -0.7927911 ]
 [ 0.43675795 -1.26224397  0.95964283  0.53281738 -0.62054885]
 [ 1.13107671  0.0891375  -0.02051028 -0.25658162  0.20175099]
 [-0.25101365 -0.83651706  0.45385709  0.36579153 -0.07140417]
 [ 0.05000226  0.02090654 -0.06632239 -0.70140683  0.91099702]
 [ 0.33570299 -1.31070844  1.04055933  0.29474841 -0.42930995]
 [ 0.1520569  -0.61360656  0.62168381 -0.96024268  0.81257321]
 [ 0.05897928 -0.03343131 -0.041293  -0.57477116  0.79554345]]
```

Now, apply the ReLU non-linear activation function, which can be implemented by invoking the `np.maximum` on `y`, such that it does element-wise comparison between each element and 0, and returns the larger value. Recall that the  $\text{ReLU}(x) = x$  if  $x > 0$ , and 0 otherwise. Quite elegant how we can implement this nonlinearity so easily as compared to other activation methods like `tanh`, `sigmoid`, etc.

Print the output `y` of the activation function. This is the final output of our unquantized fully connected layer.

```
y = np.maximum(y, 0)
print(y)
[[0.00511569 0.          0.51116489 0.57396818 0.          ]
 [1.34222938 0.          0.34179184 0.          0.39381581]]
```



```
[0.9916536  0.          0.96734553 0.53233659 0.          ]
[0.43675795 0.          0.95964283 0.53281738 0.          ]
[1.13107671 0.0891375  0.          0.          0.20175099]
[0.          0.          0.45385709 0.36579153 0.          ]
[0.05000226 0.02090654 0.          0.          0.91099702]
[0.33570299 0.          1.04055933 0.29474841 0.          ]
[0.1520569  0.          0.62168381 0.          0.81257321]
[0.05897928 0.          0.          0.          0.79554345]]
```

Let's proceed with quantizing these weights. As you can see below, we calculate  $w_{\min}$  and  $w_{\max}$  as described earlier in this section.

```
w_min = np.min(weights)
w_max = np.max(weights)
print(w_min, w_max)
-0.8586738858321132 0.7418045798990329
```

We will quantize the weights using the `quantize()` function with the parameters  $w_{\min}$ ,  $w_{\max}$  and  $b = 8$ . Let's print out the weights to verify it worked as expected. Note that the quantized weights are 0 and 255 respectively for the smallest and largest weights in the original matrix.

```
weights_quantized = quantize(weights, w_min, w_max, 8)
print(weights_quantized)
[[124  30 251  74 146]
 [  0 135 111 136 190]
 [ 96   5 228 255  16]]
```

Let's continue with the dequantization process to recover an estimate of the original weights. We will also compute the difference between the two weights using the [Root Mean Square Error](#) (RMSE), which gives an idea of the amount of information lost in the quantization process. As we see, the error is quite small.

```
weights_dequantized = dequantize(weights_quantized, w_min, w_max, 8)
weights_diff = np.sqrt(np.mean((weights_dequantized - weights) ** 2))
print(weights_diff)
0.003925407435722753
```

Now, we'll calculate the final output after the activation function and evaluate the error between the two results. Notice that the error is very small.

```
y_via_quant = np.maximum(np.matmul(inputs, weights_dequantized) + bias, 0)
y_diff = np.sqrt(np.mean((y_via_quant - y) ** 2))
print(y_diff)
0.002573583625884542
```

We are able to compress the 32-bit floating point weights tensor of shape  $[3, 5]$ , which is  $3 \times 5 \times 4 = 60$  bytes (15 elements taking up 4 bytes each) in size to a 8-bit unsigned integer compressed tensor which uses  $3 \times 5 \times 1 = 15$  bytes. It gives us a 4x savings in storage and communication costs as compared to the high precision representation. The quantized

representation size increases linearly with the increase in the value of  $b$  as shown figure 2-8 (for image sizes).

There are other works in the literature that demonstrate different variants of quantization. XNOR-Net<sup>3</sup>, Binarized Neural Networks<sup>4</sup> and others use  $b=1$ , and thus have binary weight matrices. The quantization function there is simply the  $\text{sign}(x)$  function, which is 1 if  $x$  is positive, 0 otherwise. The promise with such extreme quantization approaches is the theoretical 32x reduction in model size of larger networks like Inception without substantial quality loss. However this approach needs to be evaluated on small networks like the MobileNet.

While these approaches (and other schemes like Ternary Weight Networks<sup>5</sup>) can lead to efficient implementations of standard operations where multiplications and divisions are replaced by cheaper operations like addition and subtraction, these gains need to be evaluated in practical settings because they require support from the underlying hardware. Moreover, multiplications and divisions are cheaper at lower precisions like  $b=8$  and are well supported by the hardware technologies like the fixed-point SIMD instructions which allows data parallelism, the SSE instruction set in x86 architecture, and similar support on ARM processors as well as on specialized DSPs like the Qualcomm Hexagon.

We started out this section with two main objectives. The first one was to reduce the model size which is fulfilled using the quantization techniques. The second goal is to improve inference latency for the quantized networks. That is our next topic!

## Activation Quantization

To improve inference latency with quantized networks, we need to quantize their activations as well. It means that we will perform the math operations in the layers (including intermediate layers and the output layer) using the fixed-point quantized representation. This will save dequantization costs since the quantized weight matrices can be directly used for calculations along with the inputs.

Vanhoucke et al.<sup>6</sup> demonstrated a 3x inference speedup using a fully fixed-point model as compared to a floating-point model on an x86 CPU without sacrificing the accuracy. All the layer inputs (except the first layer) and the activations are fixed-point along with the quantized weights. The primary driver for the performance improvement was the availability of fixed-point SIMD instructions in Intel's SSE4 instruction set which can parallelize Multiply-Accumulate (MAC) operations.

---

<sup>3</sup> Rastegari, Mohammad, et al. "Xnor-net: Imagenet classification using binary convolutional neural networks." *European conference on computer vision*. Springer, Cham, 2016.

<sup>4</sup> Hubara, Itay, et al. "Binarized neural networks." *Advances in neural information processing systems* 29 (2016).

<sup>5</sup> Li, Fengfu, Bo Zhang, and Bin Liu. "Ternary weight networks." *arXiv preprint arXiv:1605.04711* (2016).

<sup>6</sup> Vanhoucke, Vincent, Andrew Senior, and Mark Z. Mao. "Improving the speed of neural networks on CPUs." (2011).

Figure 2-10 compares the accuracies and the latencies between unoptimized and quantized models. For a given latency value, the quantized variant performs with a higher accuracy in most cases. However, the best accuracy value for the unoptimized variant is slightly higher than 70% which for the quantized variant is slightly lower than 70%. And that is the precision trade off!

It's time to put our understanding of quantization into practice with a hands-on project. We will apply the learnings from weight and activation quantizations to a real world deep learning model and demonstrate the size reduction and inference efficiency improvements. The project will use the famous MNIST dataset!

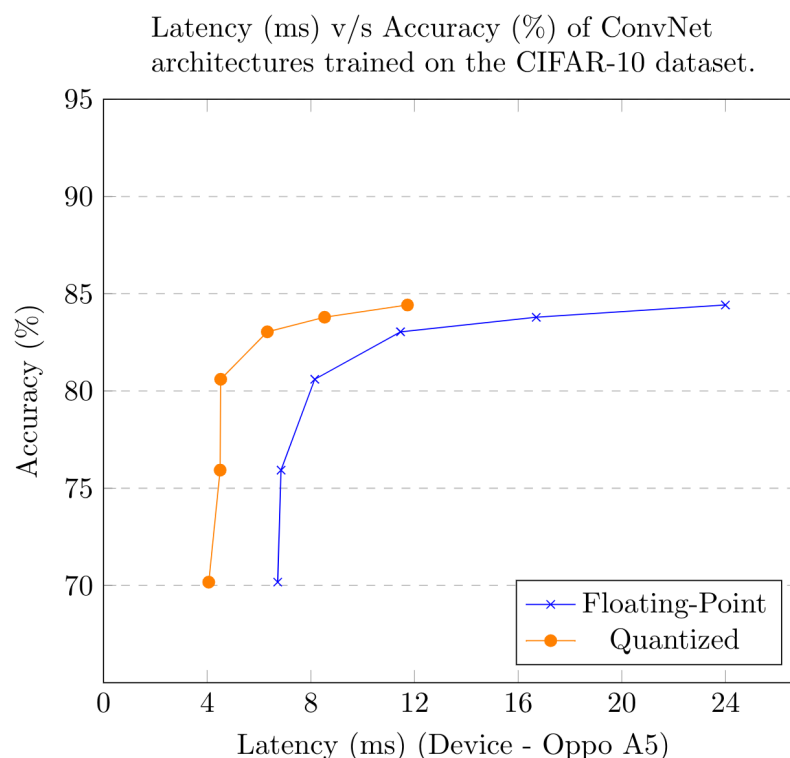


Figure 2-10: Latency v/s accuracy trade off for unoptimized representation (float) and quantized representation (8-bit) using a convolutional net trained on the CIFAR-10 dataset.

## Project: Quantizing A Deep Learning Model

We have worked through quite a bit of theory and exercises on quantization. It is time to put them into practice. [MNIST](#) (Modified NIST) handwritten digit recognition is a well-known problem in the deep learning field. We will use it to demonstrate how the quantization techniques can be applied in a practical setting by leveraging the built-in support for such technologies in the real world machine learning frameworks. We will start with the training and evaluation of a baseline (unoptimized) model. Then, we will evaluate a quantized model and compare its performance (accuracy and parameter sizes) with the baseline. Our goal with this project is to provide the readers with hands-on experience with training, evaluating

and quantizing models using the real world frameworks. Let's start with introducing the logistics!

## Logistics

In this book, we have chosen to work with Tensorflow 2.0 (TF) because it has exhaustive support for building and deploying efficient models on devices ranging from TPUs to edge devices at the time of writing. However, we encourage you to explore other frameworks like PyTorch, Apple's CoreML as well which are covered in chapter 10. If you are not familiar with the tensorflow framework, we refer you to the book *Deep Learning with Python*<sup>7</sup>. All the code examples in this book are available at the [EDL](#) github repository. The code examples for this project are available in the [IPython notebook](#) in the same repository.

You can run the code in [Google's Colab](#) environment which provides free access to CPU, GPU, and TPU resources. You can also run this locally on your machine using the Jupyter framework or with other cloud services.

## Recognizing Hand-Written Digits

With the logistics out of the way, let's solve the problem of recognizing digits on *checks or cheques* using a deep learning system. We are targeting this system to run on a low end Android device. The resource limitations are **under 50 Kb** of model size and **an upper limit of 1 millisecond** per prediction.

This sounds challenging because the human handwriting varies from person-to-person. However, there is some basic structure in handwritten digits that a neural network should be able to learn. [MNIST](#) (Modified NIST) handwritten recognition is one of the most commonly solved problems by beginners in the deep learning field. The MNIST dataset was assembled and processed by Yann LeCun et al. It is a collection of digits from 0-9 written by approximately 250 different writers including high-school students and census bureau employees. The dataset has 60,000 training examples and 10,000 test examples. Figure 2-11 shows a sample of 100 labelled images from this dataset.

Each input example is a 28x28 matrix containing values in the range [0, 255]. The task is to identify which digit class a given example belongs to. This problem can be solved with a simple deep learning model. In fact, it is one of the first tasks that [Convolutional Neural Networks](#) were used for.

---

<sup>7</sup> Deep Learning with Python by Francois Chollet

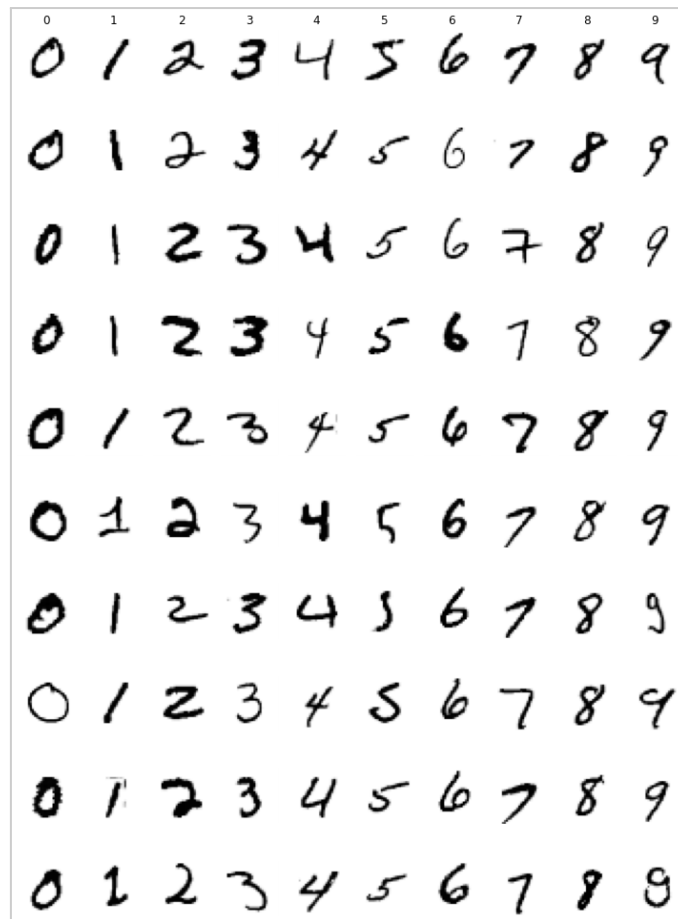


Figure 2-11: A visualization of 100 samples from the MNIST dataset.

## Loading and Processing the MNIST Dataset

Take a look at the `load_data()` function in the following code. It uses the TF's handy available MNIST dataset. We normalize the training data after it is loaded. The original data is in the range  $[0, 255]$ . To be able to process it well, we will ensure that it is in a symmetrical range between  $[-1.0, 1.0]$ . It is easy to note that this can be done by dividing by 127.5 ( $255.0 / 2$ ), and then subtracting 1.0.

We ensure that the input data is a rank-4 tensor by adding an extra dimension at the end using the `expand_dims()` method. The original data is shaped as `[number of examples, 28, 28]`. The addition of an extra dimension reshapes it as `[number of examples, 28, 28, 1]` such that each example is of shape `[28, 28, 1]`. This extra dimension makes it compatible with 2-D convolutional layers which expect one dimension for the channels. Typically with an RGB image there are 3 channels, but since there is a grayscale image there is only one channel which the dataset does not create explicitly.

The normalization and reshaping of the data is done by the `process_x()` method which is invoked for both the train and test images. Once we have loaded our data and processed it, we can do some fun stuff with it.

```
import numpy as np
```

```

import tensorflow as tf
import tensorflow.keras as keras
from keras.utils import np_utils

import tensorflow_datasets as tfds

def process_x(x):
    """Process the given tensors of images."""
    x = x.astype(np.float32)

    # The original data is in [0.0, 255.0].
    # This normalization helps in making them lie between [-1.0, 1.0].
    x /= 127.5
    x -= 1.0

    # Add one dimension for the channels.
    x = np.expand_dims(x, 3)
    return x

def load_data(ds=tf.keras.datasets.mnist):
    """Returns the processed dataset."""
    (train_images, train_labels), (test_images, test_labels) = ds.load_data()

    # Process the images for use.
    train_images = process_x(train_images)
    test_images = process_x(test_images)

    return (train_images, train_labels), (test_images, test_labels)

(train_x, train_y), (test_x, test_y) = load_data()

# You can train on the Fashion MNIST dataset, which has the exact same format
# as MNIST, and is slightly harder.
# (train_x, train_y), (test_x, test_y) =
load_data(ds=tf.keras.datasets.fashion_mnist)

```

## Creating and Compiling the Model

The `create_model()` function, described below, uses the standard tensorflow APIs to create a model. We expect an input of shape [28, 28, 1] (excluding the first batch dimension). Remember that we added an extra dimension so that convolutional layers can seamlessly work with our images. Figure 2-12 shows the detailed architecture of our model.

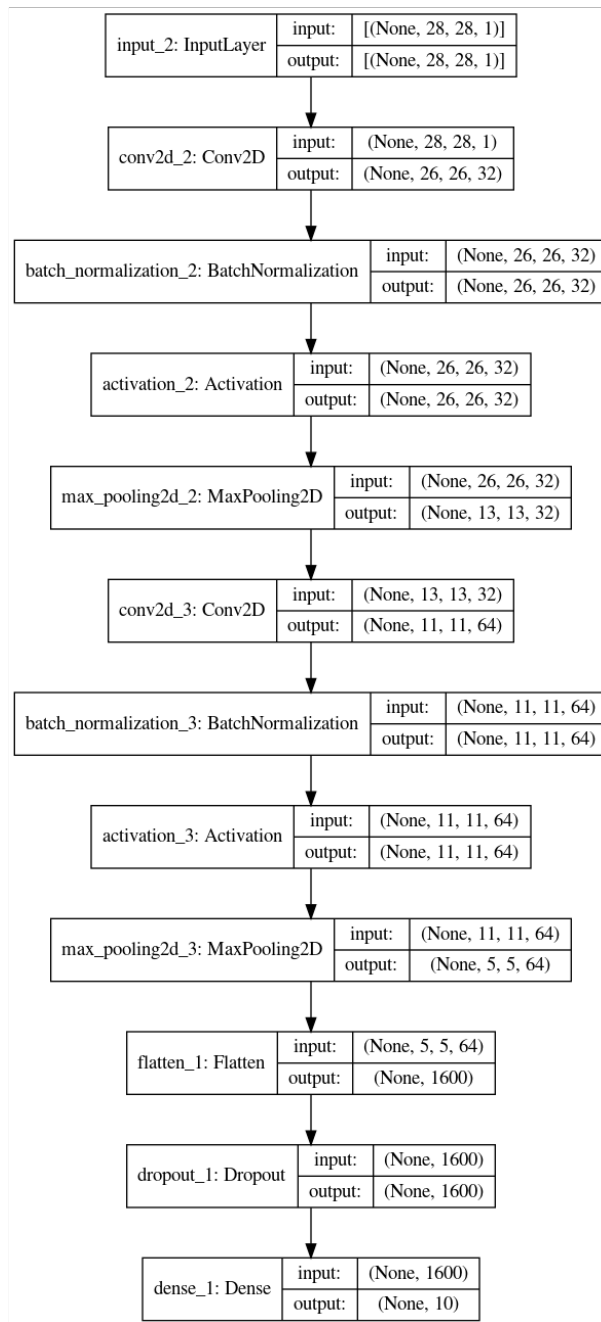


Figure 2-12: Illustration of the model that we created.

We have two convolutional layers each with filters of size 3x3. Each convolutional layer is followed by a batch-normalization layer (which rescales the output in a well-behaved range), a ReLU activation layer and a max-pooling layer to downsample the output. The first convolutional layer has 32 filters, and the second one has 64 filters. The output of the second convolution block is flattened to a rank-2 tensor (rank-1 excluding the batch dimension). A dropout layer follows right after. The final layer is a dense (fully-connected) layer of size 10 because we have 10 classes to identify.

Note that the first dimension is None, since it refers to the batch dimension which is not defined while creating the model. Our batch size could be variable (we could train with a

batch size of 16, 32, 64 and so on). The model architecture is independent of the batch size. During inference (prediction mode), the typical value for the batch size is 1 because we predict one value at a time.

The design of this model is arbitrary. You can experiment with different ideas such as stacking more convolutional layers, having a different number of filters, changing filter sizes, etc. However, be careful to adhere to certain constraints. For example, adding a dropout layer with dropout rate = 0.99 will drop most of the output and removing the non-linearity will make your model linear which will be unable to learn.

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

def create_model(dropout_rate=0.0):
    """Create a simple convolutional network."""
    inputs = keras.Input(shape=(28, 28, 1))
    x = inputs
    x = layers.Conv2D(32, (3, 3))(x)
    x = layers.BatchNormalization()(x)
    x = layers.Activation('relu')(x)
    x = layers.MaxPooling2D(pool_size=(2, 2))(x)

    x = layers.Conv2D(64, (3, 3))(x)
    x = layers.BatchNormalization()(x)
    x = layers.Activation('relu')(x)
    x = layers.MaxPooling2D(pool_size=(2, 2))(x)

    x = layers.Flatten()(x)
    x = layers.Dropout(dropout_rate)(x)
    x = layers.Dense(10)(x)
    return keras.Model(inputs=inputs, outputs=x)
```

So far, we have created a model which has stacked layers. We have also defined the input and output structure of the model. Now, let's get it ready for training. The `get_compiled_model()` function creates our model graph using the `create_model()` function. Then, it compiles the model by providing the necessary components the framework needs to train the model. This includes the loss function, the optimizer, and finally the metrics.

```
def get_compiled_model():
    """Create a compiled model (with loss fn, optimizer, metrics, etc.)"""
    model = create_model()
    model.compile(
        loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),
        optimizer=keras.optimizers.Adam(),
        metrics=tf.keras.metrics.SparseCategoricalAccuracy())
    return model
```



We use the [SparseCategoricalCrossEntropy](#) loss function as it is best suited to classify multi-class (In this case, classes are 0, 1, 2 and so on until 9) inputs. We use the sparse variant of the categorical cross entropy loss function so that we can use the index of the correct class for each example. The regular function expects one-hot labels which would require us to transform the class label 2, for example, to its one-hot representation [0 0 1 0 0 0 0 0 0].

The optimizer is the standard Adam<sup>8</sup> optimizer with the default learning rate. Feel free to tweak the learning rate and measure its impact on the training process. The metric function is SparseCategoricalAccuracy, which is the regular accuracy metric except that it knows that the labels are sparse as mentioned.

Let's print the model summary! The model has 35.2K parameters which contributes 140 KB (35.2K x 4 bytes per floating-point parameter = 140 KB) to the size of the unoptimized model. If we quantize this model, we can shave roughly 100KB from the size that will enable us to meet the size target.

```
model = get_compiled_model()
model.summary()
```

Model: "model\_1"

Layer (type)	Output Shape	Param #
input_2 (InputLayer)	[(None, 28, 28, 1)]	0
conv2d_2 (Conv2D)	(None, 26, 26, 32)	320
batch_normalization_2 (Batch Normalization)	(None, 26, 26, 32)	128
activation_2 (Activation)	(None, 26, 26, 32)	0
max_pooling2d_2 (MaxPooling2D)	(None, 13, 13, 32)	0
conv2d_3 (Conv2D)	(None, 11, 11, 64)	18496
batch_normalization_3 (Batch Normalization)	(None, 11, 11, 64)	256
activation_3 (Activation)	(None, 11, 11, 64)	0
max_pooling2d_3 (MaxPooling2D)	(None, 5, 5, 64)	0
flatten_1 (Flatten)	(None, 1600)	0
dropout_1 (Dropout)	(None, 1600)	0
dense_1 (Dense)	(None, 10)	16010
Total params: 35,210		
Trainable params: 35,018		
Non-trainable params: 192		

## Training a Unoptimized Model

We are all set to start training our model. Tensorflow provides a user-friendly API to train the model. All we need is to invoke the *fit()* method on the model object. It takes in the training

---

<sup>8</sup> Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." *arXiv preprint arXiv:1412.6980* (2014).

data, batch size, epochs, validation data and other useful parameters. The *fit()* method also prints out the training progress per epoch as shown below.

```
def train_basic_model():
    model = get_compiled_model()

    model_history = model.fit(
        train_x,
        train_y,
        batch_size=128,
        epochs=15,
        validation_data=(test_x, test_y),
        shuffle=True)

    return model, model_history.history

basic_mnist_model, basic_mnist_model_history = train_basic_model()
```

### Training Progress:

```
Epoch 1/15
469/469 [=====] - 3s 6ms/step - loss: 0.1729 -
sparse_categorical_accuracy: 0.9500 - val_loss: 0.0753 - val_sparse_categorical_accuracy:
0.9789
Epoch 2/15
469/469 [=====] - 2s 5ms/step - loss: 0.0570 -
sparse_categorical_accuracy: 0.9825 - val_loss: 0.0462 - val_sparse_categorical_accuracy:
0.9855
Epoch 3/15
469/469 [=====] - 2s 5ms/step - loss: 0.0412 -
sparse_categorical_accuracy: 0.9873 - val_loss: 0.0486 - val_sparse_categorical_accuracy:
0.9837
Epoch 4/15
469/469 [=====] - 2s 5ms/step - loss: 0.0334 -
sparse_categorical_accuracy: 0.9895 - val_loss: 0.0413 - val_sparse_categorical_accuracy:
0.9882
Epoch 5/15
469/469 [=====] - 3s 6ms/step - loss: 0.0279 -
sparse_categorical_accuracy: 0.9916 - val_loss: 0.0368 - val_sparse_categorical_accuracy:
0.9887
Epoch 6/15
469/469 [=====] - 3s 6ms/step - loss: 0.0238 -
sparse_categorical_accuracy: 0.9923 - val_loss: 0.0461 - val_sparse_categorical_accuracy:
0.9860
Epoch 7/15
469/469 [=====] - 2s 5ms/step - loss: 0.0215 -
sparse_categorical_accuracy: 0.9929 - val_loss: 0.0388 - val_sparse_categorical_accuracy:
0.9887
Epoch 8/15
469/469 [=====] - 3s 6ms/step - loss: 0.0171 -
sparse_categorical_accuracy: 0.9947 - val_loss: 0.0414 - val_sparse_categorical_accuracy:
0.9882
Epoch 9/15
469/469 [=====] - 3s 5ms/step - loss: 0.0144 -
sparse_categorical_accuracy: 0.9952 - val_loss: 0.0364 - val_sparse_categorical_accuracy:
0.9894
Epoch 10/15
469/469 [=====] - 3s 5ms/step - loss: 0.0115 -
sparse_categorical_accuracy: 0.9964 - val_loss: 0.0376 - val_sparse_categorical_accuracy:
0.9893
Epoch 11/15
469/469 [=====] - 3s 5ms/step - loss: 0.0102 -
sparse_categorical_accuracy: 0.9969 - val_loss: 0.0484 - val_sparse_categorical_accuracy:
0.9875
Epoch 12/15
469/469 [=====] - 2s 5ms/step - loss: 0.0101 -
```

```

sparse_categorical_accuracy: 0.9969 - val_loss: 0.0333 - val_sparse_categorical_accuracy:
0.9898
Epoch 13/15
469/469 [=====] - 2s 5ms/step - loss: 0.0102 -
sparse_categorical_accuracy: 0.9964 - val_loss: 0.0378 - val_sparse_categorical_accuracy:
0.9895
Epoch 14/15
469/469 [=====] - 2s 5ms/step - loss: 0.0065 -
sparse_categorical_accuracy: 0.9977 - val_loss: 0.0403 - val_sparse_categorical_accuracy:
0.9883
Epoch 15/15
469/469 [=====] - 3s 6ms/step - loss: 0.0055 -
sparse_categorical_accuracy: 0.9984 - val_loss: 0.0459 - val_sparse_categorical_accuracy:
0.9881

```

Our simple model achieved nearly 99% accuracy after 15 training epochs. The `fit()` method returns an object containing the training history which is used to plot the accuracies in figure 2-13. It shows the training and the test accuracy curves as we progressed through the training epochs.

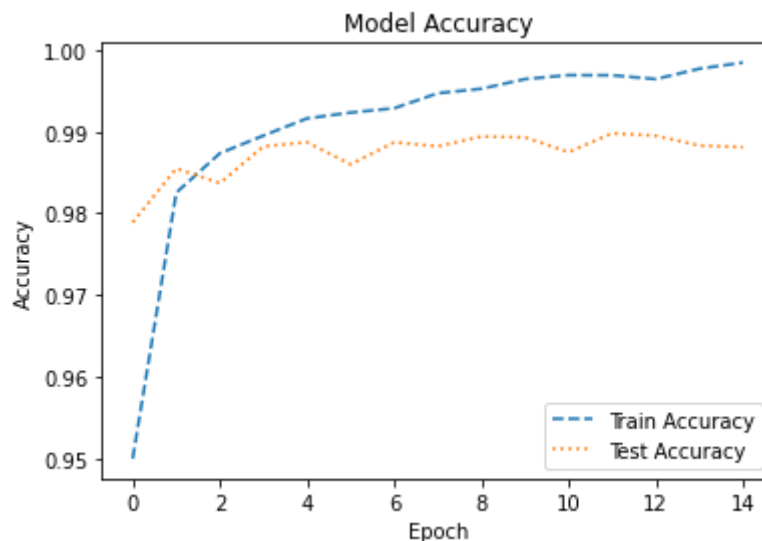


Figure 2-13: Model training with accuracy curves on train and test datasets.

## Generating an optimized TFLite Model

Now that we have got a model with good test accuracy, let's work to optimize it to fulfil the resource constraints that we initially outlined. As mentioned in Chapter 1, TFLite (Tensorflow Lite) helps to convert and deploy tensorflow models to IoT and edge devices. It is optimized for ARM based processors and supports accelerated inference using DSPs and mobile GPUs. Now, we will discuss the basics required to convert and run inference with TFLite. The details of the tensorflow ecosystem are covered in Chapter 10.

First, let's write a function, call it `tflite_model_eval()`, to evaluate a TFLite model in Python. The evaluation is a boiler-plate code. There is not much we can do to make it interesting. We are programming in the python language. Naturally, it is possible to use other languages (like Java for Android or C++ for iOS and other platforms) for inference. The authoritative guide for [TFLite inference](#) is available on the tensorflow website.

```

def tf_lite_model_eval(model_content, test_images, test_labels, quantized):
    """Evaluate the generated TFLite model."""
    # Load the TFLite model and allocate tensors.
    interpreter = tf.lite.Interpreter(model_content=model_content)
    interpreter.allocate_tensors()

    # Get input and output tensors.
    input_details = interpreter.get_input_details()
    output_details = interpreter.get_output_details()

    num_correct = 0
    num_total = 0

    for idx in range(len(test_images)):
        num_total = num_total + 1

        input_data = test_images[idx:idx+1]
        if quantized:
            # If the model is quantized, then we would have to provide the input
            # in [-127,127].
            # Rescale that data to be in [-127, 127] and then convert to int8.
            input_data = (input_data * 127).astype(np.int8)

        # Set the input tensor and invoke the interpreter.
        interpreter.set_tensor(input_details[0]['index'], input_data)
        interpreter.invoke()

        output_data = interpreter.get_tensor(output_details[0]['index'])

        # The returned output is a tensor of logits, so we find the maximum in that
        # tensor, and see if it matches the label.
        if np.argmax(output_data[0]) == test_labels[idx]:
            num_correct = num_correct + 1
    print('Accuracy:', num_correct * 1.0 / num_total)

```

The *tf\_lite\_model\_eval()* function starts by creating a tflite interpreter, which consumes the model file content. The *model\_content* variable holds the contents of the model that we created earlier. Then, it invokes *allocate\_tensors()* method on the interpreter object which allocates the space for the required tensors. The *input\_details* and *output\_details* variables are the metadata objects that describe the input and output tensors extracted from the interpreter. For a quantized model, it converts the input data to the int8 format. Earlier, we had transformed the input data to fit in the range [-1.0, 1.0]. So, we multiply the input data by 127 to map it in the range [-127.0, 127.0]. Then the *astype()* method converts it back to the int8 format. The transformed input data (for quantized models) is set as input tensor in the interpreter using the *set\_tensor()* method. We have done the hard yards. Now, we invoke the *invoke()* method to see the results.

The interpreter blocks until it finishes. Then, we use the `get_tensor()` method on the interpreter to fetch the output associated with the first output. The output is a quantized (for quantized models) logits tensor. We compute the accuracy by finding the index of the logits tensor which has the highest value. It works well with both, the quantized or the regular tensor. The logic for calculating the classification accuracy is straight-forward.

Earlier, we had referenced the `model_content` variable which is an input argument to the `tflite_model_eval()`. We stated that it holds the content of our model. What is that content? It is a transformed version of the tensorflow model that was trained in the training section. Now, let's create a combined `convert_and_eval()` function which transforms a tensorflow model to tflite format and evaluates its performance.

```
# Create the directory for storing TFLite models.
!mkdir -p 'tflite_models'

def convert_and_eval(model, model_name, quantized_export, test_dataset_x,
                    test_dataset_y):
    """Helper method to convert the given model to TFLite and eval it."""
    # Set up the converter.
    converter = tf.lite.TFLiteConverter.from_keras_model(model)

    if quantized_export:
        converter.optimizations = [tf.lite.Optimize.OPTIMIZE_FOR_SIZE,
                                  tf.lite.Optimize.OPTIMIZE_FOR_LATENCY]

    # Set up the representative dataset (using a generator function) that
    # helps improve the quality of the quantized model.
    def representative_dataset():
        for idx in range(min(len(test_dataset_x), 1000)):
            yield [test_dataset_x[idx:idx+1]]

    converter.representative_dataset = representative_dataset
    converter.target_spec.supported_ops = [tf.lite.OpsSet.TFLITE_BUILTINS_INT8]
    converter.inference_input_type = tf.int8 # or tf.uint8
    converter.inference_output_type = tf.int8 # or tf.uint8

    tflite_model_str = converter.convert()

    model_name = '{}_{}.tflite'.format(
        model_name, ('quantized' if quantized_export else 'float'))
    print('Model Name: {}, Quantized: {}'.format(model_name, quantized_export))
    print('Model Size: {:.2f} KB'.format(len(tflite_model_str) / 1024.))
    with open(os.path.join('tflite_models', model_name), 'wb') as f:
        f.write(tflite_model_str)
    # Evaluate the model.
    tflite_model_eval(tflite_model_str, test_dataset_x, test_dataset_y,
                      quantized_export)
```

As mentioned earlier, the tflite evaluation is a boiler-plate code. You can refer to the [TFLite guide](#) for more details. We start the model conversion by creating a converter object using the `from_keras_model()` method of TFLiteConverter. A call to the `convert()` method on the converter object generates a tflite model file content string. We referred to this string as `model_content` earlier. The converter object also supports weight and activation quantizations using configuration parameters.

We are almost there. We have worked out the steps to create and train a model, load and preprocess the input data and to evaluate quantized and regular models. The pieces are ready. All that is left is to put them together. Let's do that in the next section.

## Unified Training Method

We can now create a unified method that trains the model, converts it to tflite (both quantized and floating point versions), and evaluates them.

```
import os
import numpy as np

def train_model(batch_size=128, epochs=100, model_name='mnist_model'):
    model = get_compiled_model()
    model.summary()

    model_history = model.fit(
        train_x,
        train_y,
        batch_size=batch_size,
        epochs=epochs,
        validation_data=(test_x, test_y),
        shuffle=True)

    print("Running Final Evaluation")
    model.evaluate(test_x, test_y)

    # Convert and evaluate both (floating point and quantized models).
    convert_and_eval(model, model_name, False, test_x, test_y)
    convert_and_eval(model, model_name, True, test_x, test_y)

    return model, model_history.history
```

The `train_model()` function is similar to the `train_basic_model()` function. What's different is that we also invoke the `convert_and_eval()` function after the training. We invoke it twice, first with `quantized_export` set to False and then again with it set to True.

Let's train!

```
mnist_model, mnist_model_history = train_model(epochs=15)
```

Epoch 1/15

```

469/469 [=====] - 5s 9ms/step - loss: 0.1620 -
sparse_categorical_accuracy: 0.9522 - val_loss: 0.0752 - val_sparse_categorical_accuracy:
0.9771
Epoch 2/15
469/469 [=====] - 4s 9ms/step - loss: 0.0549 -
sparse_categorical_accuracy: 0.9835 - val_loss: 0.0446 - val_sparse_categorical_accuracy:
0.9849
Epoch 3/15
469/469 [=====] - 4s 9ms/step - loss: 0.0417 -
sparse_categorical_accuracy: 0.9869 - val_loss: 0.0584 - val_sparse_categorical_accuracy:
0.9806
Epoch 4/15
469/469 [=====] - 4s 9ms/step - loss: 0.0335 -
sparse_categorical_accuracy: 0.9899 - val_loss: 0.0383 - val_sparse_categorical_accuracy:
0.9871
Epoch 5/15
469/469 [=====] - 4s 9ms/step - loss: 0.0273 -
sparse_categorical_accuracy: 0.9915 - val_loss: 0.0435 - val_sparse_categorical_accuracy:
0.9867
Epoch 6/15
469/469 [=====] - 4s 9ms/step - loss: 0.0232 -
sparse_categorical_accuracy: 0.9926 - val_loss: 0.0409 - val_sparse_categorical_accuracy:
0.9879
Epoch 7/15
469/469 [=====] - 4s 9ms/step - loss: 0.0201 -
sparse_categorical_accuracy: 0.9937 - val_loss: 0.0342 - val_sparse_categorical_accuracy:
0.9893
Epoch 8/15
469/469 [=====] - 4s 9ms/step - loss: 0.0164 -
sparse_categorical_accuracy: 0.9947 - val_loss: 0.0392 - val_sparse_categorical_accuracy:
0.9890
Epoch 9/15
469/469 [=====] - 4s 9ms/step - loss: 0.0159 -
sparse_categorical_accuracy: 0.9951 - val_loss: 0.0479 - val_sparse_categorical_accuracy:
0.9864
Epoch 10/15
469/469 [=====] - 4s 9ms/step - loss: 0.0110 -
sparse_categorical_accuracy: 0.9966 - val_loss: 0.0345 - val_sparse_categorical_accuracy:
0.9906
Epoch 11/15
469/469 [=====] - 4s 9ms/step - loss: 0.0112 -
sparse_categorical_accuracy: 0.9963 - val_loss: 0.0302 - val_sparse_categorical_accuracy:
0.9904
Epoch 12/15
469/469 [=====] - 4s 9ms/step - loss: 0.0106 -
sparse_categorical_accuracy: 0.9963 - val_loss: 0.0369 - val_sparse_categorical_accuracy:
0.9892
Epoch 13/15
469/469 [=====] - 4s 9ms/step - loss: 0.0079 -
sparse_categorical_accuracy: 0.9977 - val_loss: 0.0337 - val_sparse_categorical_accuracy:
0.9905
Epoch 14/15
469/469 [=====] - 4s 9ms/step - loss: 0.0058 -
sparse_categorical_accuracy: 0.9982 - val_loss: 0.0369 - val_sparse_categorical_accuracy:
0.9893
Epoch 15/15
469/469 [=====] - 4s 9ms/step - loss: 0.0092 -
sparse_categorical_accuracy: 0.9968 - val_loss: 0.0562 - val_sparse_categorical_accuracy:
0.9862
Running Final Evaluation
313/313 [=====] - 1s 3ms/step - loss: 0.0562 -
sparse_categorical_accuracy: 0.9862
INFO:tensorflow:Assets written to: /tmp/tmp6ouf8al7/assets
Model Name: mnist_model_float.tflite, Quantized: False
Model Size: 138.71 KB
Accuracy: 0.9862
INFO:tensorflow:Assets written to: /tmp/tmpljlefd2f/assets
INFO:tensorflow:Assets written to: /tmp/tmpljlefd2f/assets
Model Name: mnist_model_quantized.tflite, Quantized: True
Model Size: 39.65 KB
Accuracy: 0.9861

```

The training output shows that the accuracy of the quantized model is just slightly less than that of the floating-point model. The quantized model is almost 1/4th the size of the floating

point model. Figure 2-14 shows the accuracy plot of the model on the training and the test datasets. We started out with a goal to create a smaller model without compromising the accuracy which we have achieved successfully. In the next chapter we will focus on the learning techniques to improve the model accuracy.

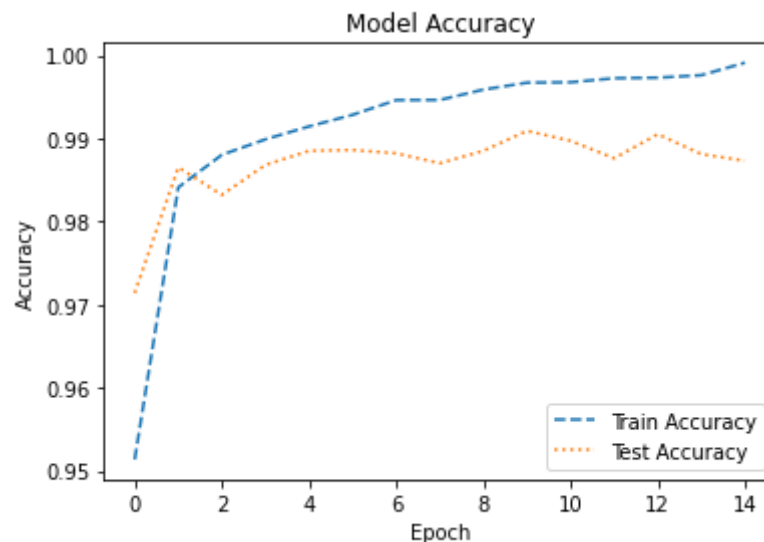


Figure 2-14: Accuracy plots of the model on train and test datasets.

## Summary

The idea of compression has been around throughout history. Nearly all of us had an opportunity to pack suitcases to move or to travel. A little bit of reorganization makes enough space for an extra pair of shoes or a couple of books to read. In the realm of the internet, videos, audios and data files are all compressed using a suitable format. It wasn't a surprise that the idea of compression crept into the deep learning field as well.

We started this chapter with a gentle introduction of compression using huffman coding and jpeg compression as examples. We talked about footprint and quality metrics as a mechanism to measure model efficiency. We learnt about quantization, a domain and model architecture agnostic compression technique that can be applied to any deep learning model. The crux of quantization is to trade off model precision for a smaller model size which results in economical storage and transmission. The mars rover example demonstrated this technique using an image of the Curiosity rover. We hope that by looking at the quality of the quantized images, the readers are able to develop an intuition on the precision trade off and compression benefits. We elaborated the idea further in the context of deep learning models by quantizing the weight matrices and demonstrated that the process is identical. Finally, we trained a model to recognize handwritten digits, quantized it and measured its performance which turned out to be almost identical to the original model. Moreover, the quantized model was 4X smaller than the original model.



Deep learning is an exciting and fast growing field which is fortunate to enjoy a large community of researchers, developers and entrepreneurs. It excites us when we come across a problem that it can solve efficiently. However, often it happens that after training a model with decent accuracy, the environmental constraints restrict the deployment of the solution for practical purposes. What we want for the reader to take away after reading this chapter is that it is not the end of the road. A giant model can be made smaller, the inference can be quicker and the memory requirements can be brought down if we are ready to make certain trade-offs. We hope that this chapter helps more deep learning models to cross the finish line. The next chapter will introduce learning techniques to improve quality metrics like accuracy and recall without impacting the model footprint. These techniques are training time techniques which are specifically useful for data scarce scenarios. Stay tuned!