
Hexadecimal Image Captcha classification using CNN

Aman Sah*

M.tech in Aerospace Engineering

Rahul Dhali

M.tech in Aerospace Engineering

Ganesh Valtule

M.Tech (Department of Aerospace Engineering)

1 Image processing, CNN, Hexadecimal Image Captcha, Logit loss, Adam optimization, parity classification
2

Abstract

3 The objective of this report is to develop a solution for classifying the parity (odd or
4 even) of hexadecimal numbers presented in CAPTCHA images. Each CAPTCHA
5 image is 500×100 pixels and contains a four-digit hexadecimal code, where the
6 characters may be rotated and rendered with different colors. The background
7 colors are light in shade, and stray lines of varying thickness and color are added
8 to enhance realism. To accomplish this task, a dataset consisting of 2000 training
9 images with corresponding labels indicating the parity of the hexadecimal numbers
10 is provided. Additionally, reference images of unrotated characters 0-9 and A-F
11 are given for reference purposes.

12 The approach to solving this problem involves developing a machine-learning
13 model capable of extracting relevant features from the CAPTCHA images and
14 making accurate predictions regarding the parity of the hexadecimal number. Tech-
15 niques such as image preprocessing, feature extraction, and classification algo-
16 rithms will be explored and evaluated to identify the most effective solution.

17 The accuracy and performance of the developed model will be assessed through
18 appropriate evaluation metrics using a validation dataset. The report will also
19 discuss any challenges encountered during the development process and potential
20 strategies for improving the model's performance.

21 Ultimately, this project aims to provide an automated solution for determining the
22 parity of hexadecimal numbers in CAPTCHA images, which can have practical
23 applications in various domains, such as web security and authentication systems.

24 1 Introduction

25 The given data set is shown below in Fig(1), representing 2000 such image captcha. Each CAPTCHA
26 image in this assignment will be 500×100 pixels in size. Each image will contain a code that is a
27 4-digit hexadecimal number. The font of all these characters would be the same, as would the font
28 size. The Latin character A-F will always be in upper case. However, each character may be rotated
29 (the degree of rotation will always be either 0° , $\pm 10^\circ$, $\pm 20^\circ$, $\pm 30^\circ$, degree and each character may
30 be rendered with a different color. The background color of each image can also change. However,
31 all background colors are light in shade. Each image also has some stray lines in the background
32 which are of varying thickness, varying color, and a shade darker than that of the background. These
33 stray lines are intended to make the CAPTCHAs are more "interesting" and realistic. Hexadecimal
34 numbers are written in base 16 i.e. there are 16 "digits" instead of the usual 10. These digits are 0, 1,
35 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F. The last six digits have values going from 10 to 15 i.e. A = 10,
36 B = 11, C = 12 etc. Thus, the hexadecimal number 10 is equal to the decimal number $1 \times 16 + 0 \times$

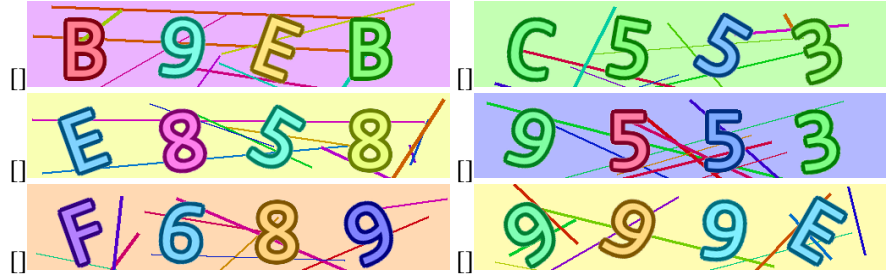


Figure 1: Random Test Set

160 = 16 and the hexadecimal number DECAF is equal to the decimal number $D \times 164 + E \times 163 + C \times 162 + A \times 161 + F \times 160 = 13 \times 164 + 14 \times 163 + 12 \times 162 + 10 \times 161 + 15 \times 160 = 912559$. Your task is to figure out whether the 4 digit hexadecimal number in the image is odd or even.

2 Methodology

Convolutional Neural Networks (CNNs) are a type of deep learning model that is widely used in image classification tasks. CNNs are inspired by the organization of animal visual cortex and designed to automatically and adaptively learn spatial hierarchies of features, from low- to high-level patterns. A CNN is typically composed of three types of layers: convolution, pooling, and fully connected layers. The convolution layer is the core building block of a CNN. It applies a set of filters to the input image and produces a set of feature maps. Each filter is a small matrix of weights that is learned during the training process. The filter slides over the input image, performing a dot product at each location, and produces a single output value. The output values are then arranged into a feature map, which highlights the presence of a particular feature in the input image..[ref.] . The pooling layer is used to reduce the spatial size of the feature maps and to make the CNN more robust to small variations in the input image. The most common pooling operation is max pooling, which takes the maximum value in a small region of the feature map and outputs it as the new value for that region. This operation reduces the size of the feature map by a factor of two.[ref.]

The fully connected layer is used to map the high-level features learned by the convolutional layers to the output classes. It takes the flattened output of the last convolutional layer and applies a set of weights to produce a vector of scores for each class. The class with the highest score is then chosen as the predicted output.

In summary, CNNs are a powerful tool for image classification tasks. They are composed of convolutional, pooling, and fully connected layers, which work together to learn spatial hierarchies of features and to map those features to the output classes. The mathematical operations involved in CNNs can be complex, but the intuition behind them is straightforward.

2.1 Image Transformation in CNN

CNN is a powerful algorithm for image processing. These algorithms are currently the best algorithms we have for the automated processing of images. Images contain data of RGB combination. Matplotlib can be used to import an image into memory from a file. The computer doesn't see an image, all it sees is an array of numbers. Color images are stored in 3-dimensional arrays. The first two dimensions correspond to the height and width of the image (the number of pixels). The last dimension corresponds to each pixel's red, green, and blue colors. Convolutional Neural Networks specialized for applications in image video recognition. CNN is mainly used in image analysis tasks like Image recognition, Object detection Segmentation.

There are three types of layers in Convolutional Neural Networks: 1) Convolutional Layer: In a typical neural network each input neuron is connected to the next hidden layer. In CNN, only a small region of the input layer neurons connect to the neuron hidden layer. 2) Pooling Layer: The pooling layer is used to reduce the dimensionality of the feature map. There will be multiple activation pooling layers inside the hidden layer of the CNN. 3) Fully-Connected layer: Fully Connected Layers form the last few layers in the network. The input to the fully connected layer is the output from the final Pooling or Convolutional Layer, which is flattened and then fed into the fully connected layer.

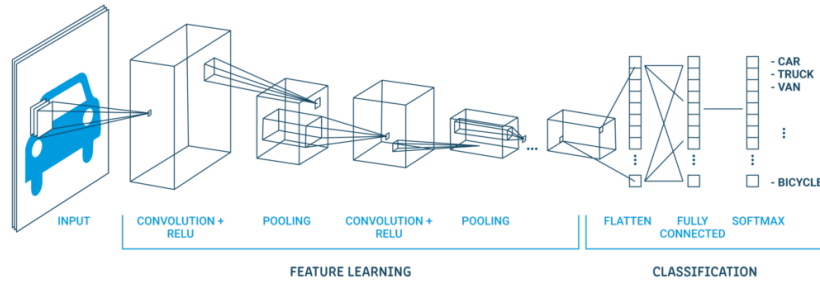


Figure 2: Example, how CNN works [ref.]

One must have come across Resnets while working with CNNs, or at least would have heard of it and we do know that ResNets perform really well on most of Computer vision tasks, but why was there even a need for an architecture like this when we already had other good performing architectures, to answer this, let us understand the drawbacks with other deep neural network architectures that were used before ResNets.

We think that the deeper the neural network, the better the performance, but when researchers experimented with deeper neural nets, it was found that adding more layers to a deep network does not always add up to its performance but rather decreases it, which was due to vanishing gradients in very deep neural networks. As a result, it was proposed that adding more layers to a deep neural network should either increase its performance or let it stay the same, but it should never decrease the performance. In order to achieve this, they came up with the concept of Skip connections/Residual connections, by use of which we can avoid loss of information flow. Let us understand what Skip connections are.

3 ResNet-18

ResNet-18 is a popular deep-learning model architecture that belongs to the family of Residual Neural Networks (ResNets). It was introduced by Kaiming He et al. in their paper "Deep Residual Learning for Image Recognition" in 2016. ResNet-18 is specifically designed for image classification tasks and has achieved excellent performance on various benchmark datasets. The key idea behind ResNet-18 is the introduction of residual blocks, which enable the network to learn residual functions instead of trying to directly fit the desired underlying mapping. This helps to alleviate the degradation problem that occurs when deeper networks suffer from increased training errors as they grow deeper. Now, let's dive deeper into the architecture and working of ResNet-18:

3.1 Basic Building Block

The fundamental building block of ResNet-18 is the residual block. Each residual block consists of two convolutional layers followed by batch normalization and a ReLU activation function. The input to the residual block is passed through the first convolutional layer, batch normalization, and activation function, and then through the second convolutional layer and batch normalization. The output of the second convolutional layer is added to the input (identity shortcut) to form the residual connection. The final output of the residual block is obtained by applying the ReLU activation function to the sum of the input and the residual connection. This formulation enables the network to learn the residual mapping, making it easier to optimize deep networks.

3.1.1 Activation function

Artificial neural networks are inspired by the biological neurons within the human body which activate under certain circumstances resulting in a related action performed by the body in response. Artificial neural nets consist of various layers of interconnected artificial neurons powered by activation functions that help in switching them ON/OFF. Like traditional machine learning algorithms, here too, there are certain values that neural nets learn in the training phase.

115 Briefly, each neuron receives a multiplied version of inputs and random weights which is then added
 116 with static bias value (unique to each neuron layer), this is then passed to an appropriate activation
 117 function which decides the final value to be given out of the neuron. There are various activation
 118 functions available as per the nature of input values. Once the output is generated from the final
 119 neural net layer, the loss function (input vs output) is calculated and backpropagation is performed
 120 where the weights are adjusted to make the loss minimum. Finding optimal values of weights is what
 121 the overall operation is focusing around.
 122 As mentioned above, activation functions give out the final value given out from a neuron, but what is
 123 an activation function and why do we need it?

124 So, an activation function is basically just a simple function that transforms its inputs into outputs
 125 that have a certain range. There are various types of activation functions that perform this task in a
 126 different manner. For example, the sigmoid activation function takes input and maps the resulting
 127 values in between 0 to 1.

128 One of the reasons that this function is added to an artificial neural network in order to help the
 129 network learn complex patterns in the data. These functions introduce nonlinear real-world
 130 properties to artificial neural networks. Basically, in a simple neural network, x is defined as inputs,
 131 w weights, and we pass $f(x)$ which is the value passed to the output of the network. This will then be
 132 the final output or the input of another layer.

133 If the activation function is not applied, the output signal becomes a simple linear function. A neural
 134 network without activation function will act as a linear regression with limited learning power. But
 135 we also want our neural network to learn non-linear states as we give it complex real-world
 136 information such as images, video, text, and sound.

137 3.2 ReLU Activation Function

138 ReLU stands for rectified linear activation unit and is considered one of the few milestones in the
 139 deep learning revolution. It is simple yet really better than its predecessor activation functions such
 140 as sigmoid or tanh. ReLU activation function formula, First, let us define a ReLU function. ReLU
 141 function is its derivative both are monotonic. The function returns 0 if it receives any negative input,
 142 but for any positive value x , it returns that value back. Thus it gives an output that has a range from 0
 143 to infinity. Now let us give some inputs to the ReLU activation function and see how it transforms.

$$144 \quad f(x) = \max(0, x)$$

```
145 def ReLU(x):
146     if x > 0:
147         return x
148     else:
149         return 0

150 input_series = [x for x in range(-19, 19)]
151 calculate outputs for our inputs
152 output_series = [ReLU(x) for x in input_series]
153
```

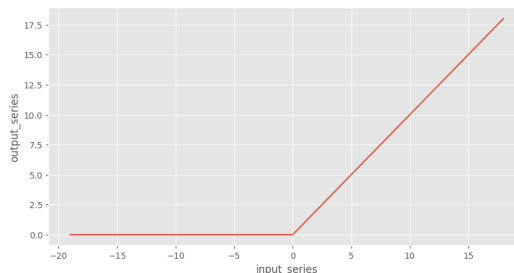


Figure 3: Relu function example graph

ReLU is used as a default activation function and nowadays and it is the most commonly used activation function in neural networks, especially in CNN. Next, we will give input as a list of numbers from -19 to +19. the ReLU function is simple and it consists of no heavy computation as there is no complicated math. The model can, therefore, take less time to train or run. One more important property that we consider the advantage of using the ReLU activation function is sparsity. Usually, a matrix in which most entries are 0 is called a sparse matrix and similarly, we desire a property like this in our neural networks where some of the weights are zero. Sparsity results in concise models that often have better predictive power and less overfitting/noise. In a sparse network, it's more likely that neurons are actually processing meaningful aspects of the problem. For example, in a model detecting human faces in images, there may be a neuron that can identify ears, which obviously shouldn't be activated if the image is not of a face and is a ship or mountain. Since ReLU gives output zero for all negative inputs, it's likely for any given unit to not activate at all which causes the network to be sparse. Now let us see how the ReLU activation function is better than previously famous activation functions such as sigmoid and tanh.

3.3 ResNet-18 Architecture

ResNet-18 consists of several stacked residual blocks, forming the overall architecture of the model. The initial layer is a standard convolutional layer with a large kernel size (7x7) and a stride of 2, followed by batch normalization and a ReLU activation function. This is followed by a max-pooling layer. Then, four stages of residual blocks are stacked together, each with a different number of residual blocks. The number of filters in the convolutional layers gradually increases from stage to stage, while the spatial dimensions are reduced through stridden convolutions. Finally, an adaptive average pooling layer is applied to reduce the spatial dimensions to a fixed size, and a fully connected layer is used for the final classification.

3.4 Shortcut Connections

The identity shortcut connections in ResNet-18 are crucial for addressing the degradation problem. By adding the input to the output of the residual block, the network can learn to make small incremental changes to the identity mapping, rather than trying to learn the entire mapping from scratch. These shortcuts also help propagate gradients more effectively during backpropagation, allowing for better optimization of deeper networks

3.5 Training and Optimization

ResNet-18 is typically trained using a variant of stochastic gradient descent (SGD) called mini-batch SGD. The loss function used depends on the specific task, but commonly used loss functions include cross-entropy loss for classification tasks. During training, the model's parameters are updated to minimize the loss between the predicted outputs and the ground truth labels using backpropagation and gradient descent optimization algorithms

3.6 Transfer Learning: ResNet-18

like other deep learning models, can benefit from transfer learning. Transfer learning involves leveraging the pre-trained weights of a model that was trained on a large dataset (such as ImageNet) and fine-tuning it on a smaller, task-specific dataset. By using pre-trained weights as initializations, the model can start with better representations and learn task-specific features more efficiently. In summary, ResNet-18 is a deep convolutional neural network architecture that addresses the degradation problem associated with training deep networks. By introducing residual connections, the model can learn residual mappings and optimize deeper networks more effectively. This architecture has achieved state-of-the-art results in image classification tasks and serves as a foundation for more advanced ResNet variants.

4 skip connections

A Skip/Residual connection takes the activations from an (n-1) convolution layer and adds it to the convolution output of (n+1) layer and then applies ReLU on this sum, thus Skipping the n layer. Fig.3

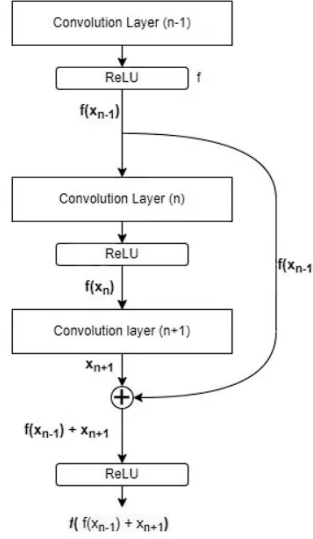


Figure 4: Skip Connections[ref.]

explains how a skip connection works. (Here I am using $f(x)$ to denote Relu applied on x where x is the output after applying the Convolution operation. But how does this even help, simply put, if the n layer is not learning anything even then we won't lose any information, because at $(n+1)$ we are using the output of the $(n-1)$ layer as well when we move forward and then applying activation on this sum. Thus we are enabling the network to skip one ReLU activation in between if it does not provide any useful information or provides no information at all i.e., 0, and the network will be using the previous information, thus maintaining consistent performance. If anyways both layers are providing significant information, thus having previous information will anyways boost the performance.

5 Loss Function And Optimizations

5.1 Loss Function

Our code uses the `nn.BCEWithLogitsLoss` class as the loss function. This loss function is suitable for binary classification problems where the model predicts a single scalar value for each input sample. It combines a sigmoid activation function and binary cross-entropy loss. `nn.BCEWithLogitsLoss` is an implementation of the binary cross-entropy loss, which is commonly used when dealing with binary classification tasks. It operates on logits, which are the output of the final layer of the model before applying a sigmoid activation function.

5.2 BCEWithLogitsLoss

This loss combines a Sigmoid layer and the BCELoss in one single class. This version is more numerically stable than using a plain Sigmoid followed by a BCELoss as, by combining the operations into one layer, we take advantage of the log-sum-exp trick for numerical stability. The unreduced (i.e. with reduction set to 'none') loss can be described as where N is the batch size. If a reduction is 'none' (default 'mean'), then

$$\ell(x, y) = L = \begin{bmatrix} l_1 \\ \vdots \\ l_N \end{bmatrix}^T, \quad l_n = -w_n [y_n \cdot \log \sigma(x_n) + (1 - y_n) \cdot \log(1 - \sigma(x_n))]$$

$$\ell(x, y) = \begin{cases} \text{mean}(L), & \text{if reduction = 'mean'}; \\ \text{sum}(L), & \text{if reduction = 'sum'}. \end{cases}$$

225 This is used for measuring the error of a reconstruction in for example an auto-encoder. Note that the
 226 targets $t[i]$ should be numbers between 0 and 1. It's possible to trade off recall and precision by
 227 adding weights to positive examples. In the case of multi-label classification, the loss can be
 228 described as:

$$\ell_c(x, y) = L_c = \begin{bmatrix} l_{1,c} \\ \vdots \\ l_{N,c} \end{bmatrix}^\top, \quad l_{n,c} = -w_{n,c} [p_c y_{n,c} \cdot \log \sigma(x_{n,c}) + (1 - y_{n,c}) \cdot \log(1 - \sigma(x_{n,c}))]$$

229 The loss function for binary classification with the class number c is defined as:

230 where C is the class number, N is the number of samples in the batch, and P_c is the weight of the
 231 positive answer for class C . $P_c > 1$ increases the recall, while $P_c < 1$ increases the precision. For
 232 example, if a dataset contains 100 positive and 300 negative examples of a single class, then
 233 *Pos_weight* for the class should be equal to $\frac{300}{100} = 3$. The loss would act as if the dataset contains
 234 $3 \times 100 = 300$ positive examples.

235 5.3 Binary Cross Entropy (BCE)

236 Binary Cross Entropy (BCE) is a common loss function used in binary classification tasks. It
 237 measures the dissimilarity between the predicted probability distribution and the true binary labels.
 238 In binary classification, we have two classes: positive and negative. The BCE loss is designed to
 239 quantify the error or discrepancy between the predicted probabilities and the true binary labels.

240 In many cases, the output of a model's final layer is transformed using the sigmoid activation function
 241 to obtain probabilities that range between 0 and 1. The sigmoid function maps the raw outputs, also
 242 known as logits, to the probability of the positive class. However, there are situations where the
 243 sigmoid activation is not applied, and the model's output consists of unnormalized scores or logits.

244 The BCElogit loss is specifically designed to handle this scenario. It applies the BCE loss directly to
 245 the logits, without the need for a sigmoid transformation. The purpose of the BCElogit loss is to
 246 penalize incorrect predictions and encourage the model to output higher logits for positive examples
 247 and lower logits for negative examples. Now, let's break down the formula for the BCElogit loss step
 248 by step:

249 $\max(\text{logit}, 0)$: This term ensures that the loss is non-negative. It is used to handle cases where the
 250 predicted logit is negative. If the predicted logit is greater than or equal to zero, the $\max(\text{logit}, 0)$
 251 term evaluates to logit itself. However, if the predicted logit is negative, the term becomes zero. This
 252 component ensures that incorrect predictions are penalized.

253 $\text{logit} * \text{label}$: This term is subtracted from the previous component. It represents the element-wise
 254 multiplication between the predicted logit and the true binary label. When the true label is 1, this
 255 term encourages higher logits for positive examples. Conversely, when the true label is 0, the term
 256 has no effect on the loss since multiplying by 0 results in 0.

257 $\log(1 + \exp(-\text{abs}(\text{logit})))$: This term is added to the previous components. It helps to smooth the loss
 258 function and stabilize the optimization process. The term involves taking the absolute value of the
 259 predicted logit, negating it, and passing it through the exponential function. The logarithm is applied
 260 to ensure numerical stability. This component contributes to the loss by increasing its value when the
 261 predicted logit deviates significantly from zero. By combining these components, the BCElogit loss
 262 encourages the model to output higher logits for positive examples and lower logits for negative
 263 examples. It penalizes incorrect predictions and helps train the model to make better binary
 264 classification decisions.

265 It's worth mentioning that different frameworks and libraries might implement the BCElogit loss
 266 with slight variations for numerical stability and efficiency. The formula provided here represents the
 267 core idea of the BCElogit loss, but specific implementations may differ.

$$\text{BCElogit loss} = \max(\text{logit}, 0) - \text{logit} \cdot \text{label} + \log(1 + \exp(-|\text{logit}|)) \quad (1)$$

268 5.4 Adam optimizer

269 5.5 Adam optimizer with BCElogit loss

270 The Adam optimizer is a popular optimization algorithm used for training neural networks. It
271 combines the benefits of AdaGrad and RMSprop. The BCElogit loss is the binary cross-entropy loss
272 applied directly to the logits, without a sigmoid transformation.

273 5.6 Mathematical Formulation

274 1. Initialize the Adam optimizer variables:

```
275         m = 0 (initialization of first-moment estimate)
276         v = 0 (initialization of second-moment estimate)
277         beta1 = 0.9 (exponential decay rate for the first-moment estimate)
278         beta2 = 0.999 (exponential decay rate for the second moment estimate)
279         epsilon = 1e-8 (small constant to prevent division by zero)
280
```

281 2. Compute the gradient of the BCElogit loss with respect to the model's parameters:

```
282         gradient = compute_gradient(loss, parameters)
283
```

284 3. Update the first-moment estimate:

```
285         m = beta1 * m + (1 - beta1) * gradient
286
```

287 4. Update the second-moment estimate:

```
288         v = beta2 * v + (1 - beta2) * (gradient ** 2)
289
```

290 5. Compute the bias-corrected first-moment estimate:

```
291         m_hat = m / (1 - beta1 ** t) (t is the current iteration or time step)
292
```

293 6. Compute the bias-corrected second-moment estimate:

```
294         v_hat = v / (1 - beta2 ** t)
295
```

296 7. Update the model's parameters using the Adam update rule:

```
297         parameters = parameters - (learning_rate * m_hat) / (sqrt(v_hat) + epsilon)
298
```

299 5.7 Background Code Implementation

300 Here's an example of how you can implement the Adam optimizer with the BCElogit loss in Python
301 using PyTorch:

```
302 import torch
303 import torch.nn as nn
304 import torch.optim as optim
305
306 # Define your model
307 model = YourModel()
308
309 # Define the BCElogit loss
310 loss_function = nn.BCEWithLogitsLoss()
311
312 # Define the Adam optimizer
313 optimizer = optim.Adam(model.parameters(), lr=0.001, betas=(0.9, 0.999), eps=1e-8)
```



```

314
315 # Training loop
316 for epoch in range(num_epochs):
317     # Clear gradients
318     optimizer.zero_grad()
319
320     # Forward pass
321     logits = model(inputs)
322     loss = loss_function(logits, labels)
323
324     # Backward pass
325     loss.backward()
326
327     # Update model parameters
328     optimizer.step()

```

329 In this code snippet, replace `YourModel` with the actual model class or instance you are using. Make
330 sure to import the necessary modules, including `torch`, `torch.nn`, and `torch.optim`.

331 Next, define the BCElogit loss function using `nn.BCEWithLogitsLoss()`. This loss function is
332 suitable for binary classification tasks where the model's final output is not passed through a sigmoid
333 activation function.

334 Then, define the Adam optimizer using `optim.Adam()` and provide the model's parameters along
335 with other required arguments like the learning rate (`lr`), beta values (`betas`), and epsilon (`eps`).
336 Adjust these hyperparameters based on your specific requirements.

337 Inside the training loop, start by clearing the gradients of the optimizer using
338 `optimizer.zero_grad()` to prevent them from accumulating between iterations.

339 Proceed with the forward pass through the model to obtain the logits using `model(inputs)`. Pass
340 the logits and the corresponding labels to the BCElogit loss function using
341 `loss_function(logits, labels)` to compute the loss.

342 Perform the backward pass by calling `loss.backward()` to compute the gradients of the loss with
343 respect to the model's parameters.

344 Finally, update the model's parameters using the `optimizer.step()` method, which applies the
345 Adam optimization algorithm to update the model's parameters based on the computed gradients and
346 the selected learning rate.

347 You can customize this code snippet by replacing `inputs` and `labels` with your own input data and
348 labels, respectively. Additionally, modify the number of epochs (`num_epochs`) to control the
349 duration of training.

350 Remember to adapt the code to your specific use case and framework, ensuring that the necessary
351 packages are imported and the model, loss function, and optimizer are appropriately defined.

352 5.8 Our Model's Code

```

353
354 from torch import nn
355
356
357 class Block(nn.Module):
358
359     def __init__(self, in_channels, out_channels, identity_downsample=None, stride=1):
360         super(Block, self).__init__()
361         self.conv1 = nn.Conv2d(in_channels, out_channels, kernel_size=3, stride=stride, padding=1)
362         self.bn1 = nn.BatchNorm2d(out_channels)
363         self.conv2 = nn.Conv2d(out_channels, out_channels, kernel_size=3, stride=1, padding=1)
364         self.bn2 = nn.BatchNorm2d(out_channels)
365         self.relu = nn.ReLU()

```

```

366         self.identity_downsample = identity_downsample
367
368     def forward(self, x):
369         identity = x
370         x = self.conv1(x)
371         x = self.bn1(x)
372         x = self.relu(x)
373         x = self.conv2(x)
374         x = self.bn2(x)
375         if self.identity_downsample is not None:
376             identity = self.identity_downsample(identity)
377         x += identity
378         x = self.relu(x)
379         return x
380
381
382 class ResNet_18(nn.Module):
383
384     def __init__(self, image_channels, num_classes):
385
386         super(ResNet_18, self).__init__()
387         self.in_channels = 64
388         self.conv1 = nn.Conv2d(image_channels, 64, kernel_size=7,
389                                 stride=2, padding=3)
390         self.bn1 = nn.BatchNorm2d(64)
391         self.relu = nn.ReLU()
392         self.maxpool = nn.MaxPool2d(kernel_size=3,
393                                       stride=2, padding=1)
394
395         #resnet layers
396         self.layer1 = self.__make_layer(64, 64, stride=1)
397         self.layer2 = self.__make_layer(64, 128, stride=2)
398         self.layer3 = self.__make_layer(128, 256, stride=2)
399         self.layer4 = self.__make_layer(256, 512, stride=2)
400
401         self.avgpool = nn.AdaptiveAvgPool2d((1, 1))
402         self.fc = nn.Linear(512, num_classes)
403
404     def __make_layer(self, in_channels, out_channels, stride):
405
406         identity_downsample = None
407         if stride != 1:
408             identity_downsample = self.identity_downsample(in_channels, out_channels)
409
410         return nn.Sequential(
411             Block(in_channels, out_channels,
412                  identity_downsample=identity_downsample, stride=stride),
413             Block(out_channels, out_channels)
414         )
415
416     def forward(self, x):
417
418         x = self.conv1(x)
419         x = self.bn1(x)
420         x = self.relu(x)
421         x = self.maxpool(x)
422
423         x = self.layer1(x)
424         x = self.layer2(x)

```

```

425         x = self.layer3(x)
426         x = self.layer4(x)
427
428         x = self.avgpool(x)
429         x = x.view(x.shape[0], -1)
430         x = self.fc(x)
431         return x
432
433     def identity_downsample(self, in_channels, out_channels):
434
435         return nn.Sequential(
436             nn.Conv2d(in_channels, out_channels,
437                       kernel_size=3, stride=2, padding=1),
438             nn.BatchNorm2d(out_channels)
439         )
440

```

441 The below table explains the model summary.

442 5.9 Model Summary

Table 1: Summary of the Model

Layer (type)	Output Shape	Param #
Conv2d-1	[32, 64, 50, 250]	12,608
BatchNorm2d-2	[32, 64, 50, 250]	128
ReLU-3	[32, 64, 50, 250]	0
MaxPool2d-4	[32, 64, 25, 125]	0
Conv2d-5	[32, 64, 25, 125]	36,928
BatchNorm2d-6	[32, 64, 25, 125]	128
ReLU-7	[32, 64, 25, 125]	0
...
AdaptiveAvgPool2d-67	[32, 512, 1, 1]	0
Linear-68	[32, 1]	513

443 Total params: 12,561,217

444 Trainable params: 12,561,217

445 Non-trainable params: 0

446 Input size (MB): 24.41

447 Forward/backward pass size (MB): 2079.98

448 Params size (MB): 47.92

449 Estimated Total Size (MB): 2152.32

451 6 Results and Conclusion

452 The training process is initiated by printing "Training" and calling the `train_step()` function. This
453 function likely performs the training step for the model, including forward and backward passes,
454 updating weights using the optimizer, calculating the loss, and computing the accuracy. The testing
455 process is initiated by printing "Testing" and calling the `test_step()` function. This function likely
456 performs the evaluation step for the model, calculating the loss and accuracy on the test dataset. The
457 loop continues for each epoch, printing the progress and updating the model. After the loop ends, the
458 timer is stopped using `time.perf_counter()` and the elapsed time is stored in the `model_4_time`
459 variable. Finally, the execution time of the model is printed.

460 6.1 Training and Testing Results

461 The code includes a training loop that runs for a specified number of epochs, performing training and
462 testing at each epoch. The training results, such as loss and accuracy, are printed during each training

iteration, while the testing results are printed after each epoch. Upon completion of the first few epochs, the accuracy of the model is found to be around 0.5, but after 7 epochs, it is found to be nearly 1.0.

6.2 Loss and Accuracy

The graph of training and testing loss with respect to the number of epochs can be seen below:

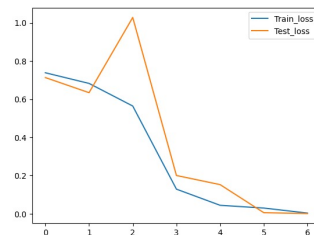


Figure 5: Training and Testing Loss with epochs

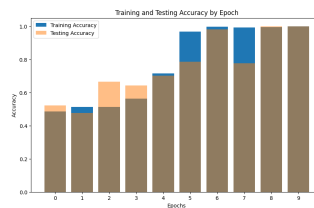


Figure 6: Training and Testing Accuracy with epochs

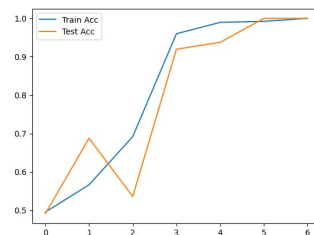


Figure 7: Training and Testing Accuracy with epochs

The graph of training and testing accuracy with respect to the number of epochs can be seen below: The accuracy graph is shown below, from here we see that the model used and the hyperparameter tuning we did are acceptable and we reached o 100 percent accuracy on the secret test data set.

7 References

1. Gentle Dive into Math Behind Convolutional Neural Networks[here]
2. Dive into CNN[here]
3. Residual Networks (ResNet) and ResNeXt[here]
4. Machine Learning Tutorial For Beginners[here]
5. ResNet — Understand and Implement from scratch.here
6. BCEWITHLOGITSLoss .here. source code
8. Gentle Introduction to the Adam Optimization Algorithm for Deep Learning.here
7. Mathematics for Machine Learning .here