

## 2.2. Machine learning models

Throughout the project, ~~different~~<sup>Various</sup> machine-learning algorithms will be used to train and evaluate the effect of using contrastive learning. This section provides a summary of these algorithms.

### 2.2.1. Logistic regression

Logistic regression is an established statistical model that is used to predict target values based on historical data. Logistic regression estimates the target values  $\hat{y}$  by applying the sigmoid function  $\sigma(\cdot)$  to the weighted input data with a constant factor  $b$ :

$$\hat{y} = \sigma(z) \quad (2.1)$$

a linear transformation of the

this is included in a "linear transformation"

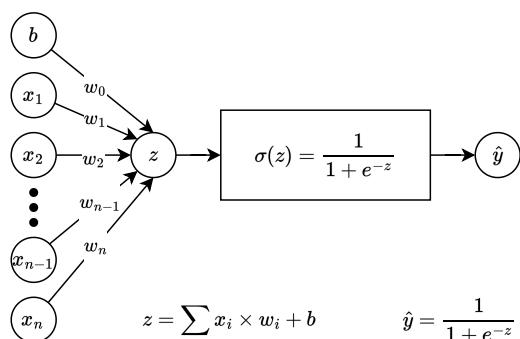
where  $z$  is calculated as:

$$z = \sum_i x_i \times w_i + b \quad (2.2)$$

and the sigmoid function is defined as

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (2.3)$$

and maps values between  $(-\infty, \infty)$  to a value between  $(0, 1)$ . The model takes in the input data  $x_i$  and multiplies each <sup>component of the</sup> input vector with a weight  $w_i$ , where the weights are optimised during training to maximise the likelihood criterion. The result is passed to the sigmoid/standard logistic function to map the results between 0 and 1. The output can be seen as the probability that the particular input variable belongs to a certain class, i.e.  $\hat{y} = P(y = 1 | \mathbf{x}, \mathbf{w})$  where  $y = 1$  if it belongs to the specified class, and 0 otherwise. Logistic regression can be seen as a single <sup>layer</sup> multi-layer perceptron with a sigmoid activation function as shown in Figure 2.1.



**Figure 2.1:** Logistic regression is shown as a multi-layer perception with input  $x$ , a sigmoid activation function  $\sigma(\cdot)$  and output  $\hat{y}$ .

*Interpret*

1  
you do not have vector output in your figure.  
Can be viewed as a single node of a

## 2.2.2. Multi-layer perceptron

A neuron is a biological structure in the brain that ~~transmits information throughout the model~~ processes and ~~information throughout the model~~ nerve signals. The artificial neuron, also referred to as a perceptron, is based on the biological neuron as shown in Figure 2.2. It receives input  $\mathbf{x}$  from a source, adds a level of importance  $\mathbf{w}$  to each input signal and generates outputs  $\mathbf{y}$ .

Suggest use lots of this figure with Fig 2.1 adapted to look more like lots of this figure but with one output  $y$   
Can click in caption.

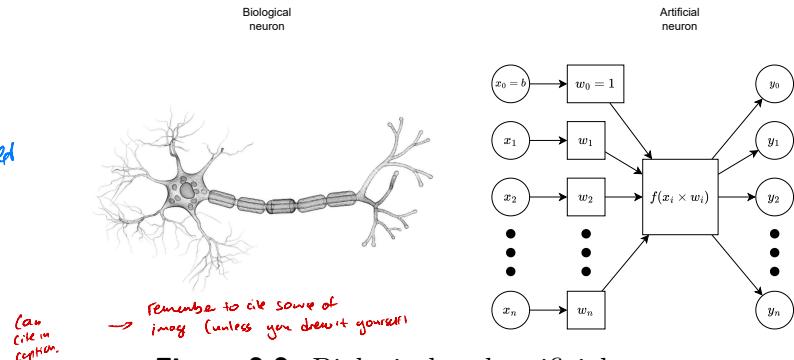


Figure 2.2: Biological and artificial neuron.

Here you have vector output, is this multinomial logistic regression?

If so, dim of input and output generally diff ( $m, n$ ) and there are  $m \times n$  weights and  $f$  becomes a softmax update?

Perhaps the several outputs are meant to represent the nerve endings of the axon in the left-hand diagram? If so note that the perceptron abstraction of the biological neuron has only one output

A multi-layer perceptron (MLP) consists of ~~multiple~~ many of these neurons/perceptrons as shown in Figure 2.3. It generally consists of input  $\mathbf{x}$ ,  $M$  hidden layers, ~~an activation function  $\sigma(\cdot)$~~  and the final output  $\hat{\mathbf{y}}$ , where the final output can be either a single value for binary classification or multiple values for multinomial classification.

Arranged in layers  
? each neuron has one output

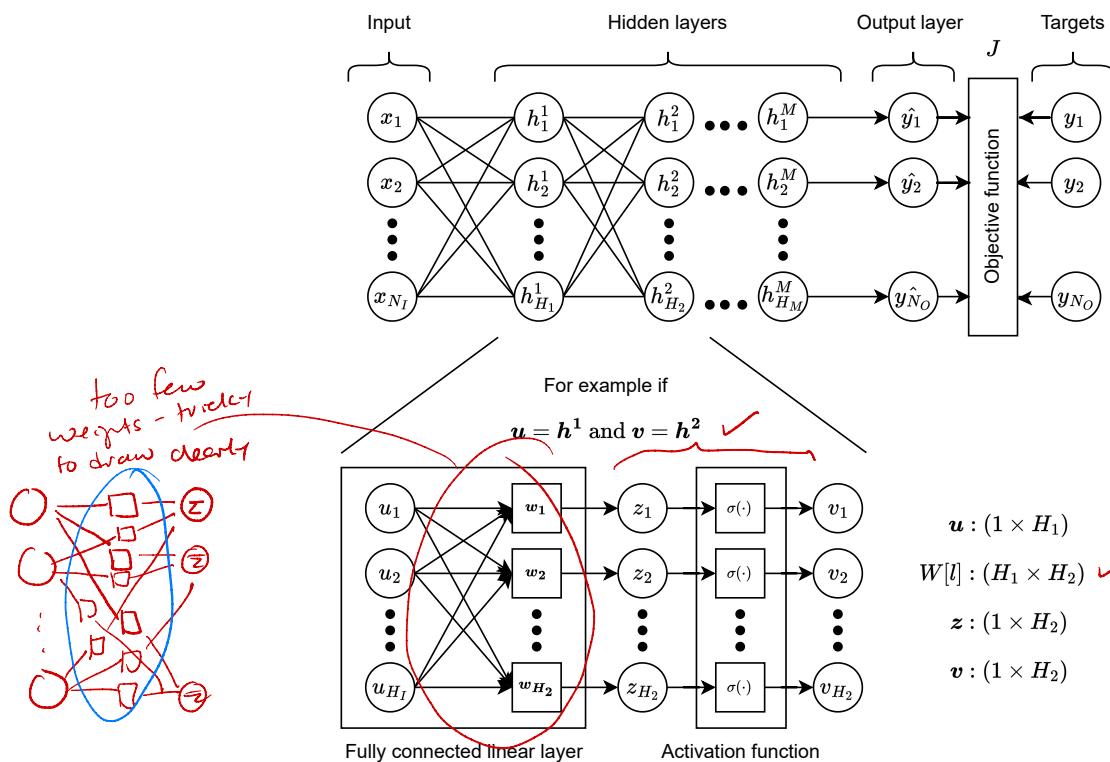


Figure 2.3: Multi-layer perceptron with input layer  $x$  with  $N_I$  inputs,  $M$  hidden layers each with  $H_1, H_2, \dots, H_M$  nodes, and predictions  $\hat{\mathbf{y}}$ .

Outputs?

I know this is going to seem impossibly pedantic, but I dislike "predictive modelling". 2.2. Machine learning models

7

time series. How about plain old "output" or perhaps "estimate" (since  $\hat{y}$  estimates  $y$ ) followed by a backward pass.

capital matrix

matrix

resulting vector

A multi-layer perceptron is trained using a forward pass. First, the input data and the weights are initialised. For each layer in the MLP, the data  $\mathbf{u}$  is multiplied by the weight  $\mathbf{w}$ . A non-linear transformation  $\sigma(\cdot)$ , referred to as the activation function, is applied to the result, to generate the output of each layer  $\mathbf{v}$ . The output of the current layer is then used as input to the next layer and the output of the last layer is the predicted values generated by the model  $\hat{\mathbf{y}}$ . The predictions are compared to the true labels or targets using a cost function/objective function. For example, if the squared error is used as an objective function, then

$$J = \frac{1}{2}(\mathbf{y} - \hat{\mathbf{y}})^T(\mathbf{y} - \hat{\mathbf{y}}) \quad (2.4)$$

To optimise the weights of the multi-layer perceptron, the objective function is minimised - allowing for the best estimation of the true labels. This occurs during a process known as back-propagation. After the model has passed through the entire forward pass, the derivative of the objective function  $J$  relative to the weights for each layer  $W[l]$  is calculated:

$$\frac{\partial J}{\partial W[l]} = \left[ \frac{\partial J}{\partial \mathbf{w}_1[l]}, \frac{\partial J}{\partial \mathbf{w}_2[l]}, \dots, \frac{\partial J}{\partial \mathbf{w}_{H_l}[l]} \right] \quad (2.5)$$

for node  $i$  in layer  $l$ :

$$\frac{\partial J}{\partial \mathbf{w}_i[l]} = \frac{\partial z[l]}{\partial \mathbf{w}_i[l]} \times \frac{\partial v[l]}{\partial z[l]} \times \frac{\partial J}{\partial v[l]} \quad (2.6)$$

These derivatives are stored and used during gradient descent to update the weights. The pseudo-code for the forward pass is provided in the Appendix, C.1. Commonly used activation functions include linear, sigmoid, softmax, tanh and ReLU. A summary of each of these activation functions is provided in the Appendix, A.1.

are the outputs of the network as a whole.

just choose one refine

$J$  is minimised during gradient descent.

The derivatives needed for gradient descent.

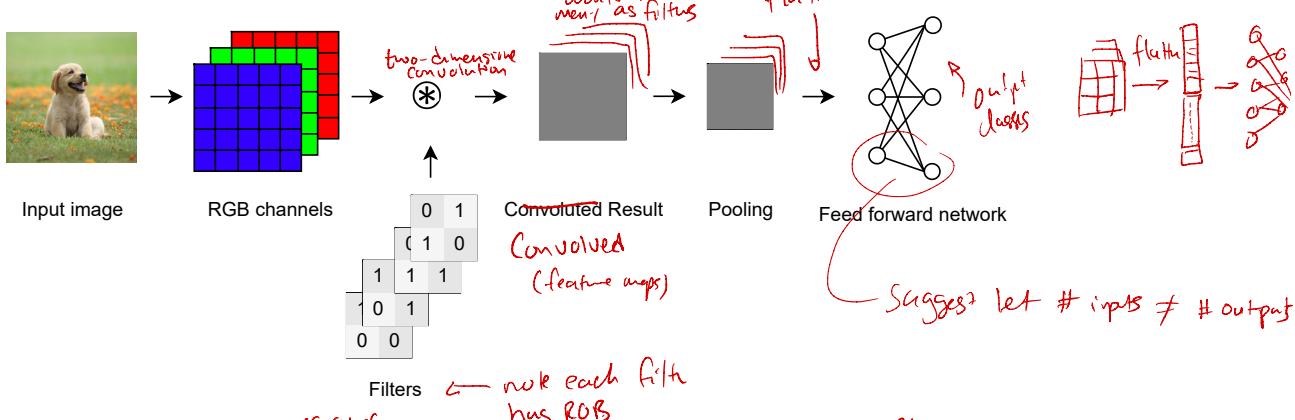
computed during backpropagation.

### 2.2.3. Convolutional neural network

A convolutional neural network (CNN) consists of one or more convolutional layers, filters, pooling layers and fully connected layers. It is a deep learning technique that takes advantage of the concept of filters to perform feature extraction on images. A simple convolutional network can be seen in Figure 2.4 to show how an image is processed, convolved, pooled and then used as input to a normal feedforward network.

*Convolved*

*presented*



**Figure 2.4:** A convolutional neural network receives input in the form of *images*. The input image is *usually* composed of three channels: red, green and blue (RGB). These channels are then convolved with the filters to perform feature detection. Pooling is used to reduce the dimension of the resulting feature maps. Finally, the output of the convolutional network is fed to a fully connected feed-forward network.

A convolutional layer uses one or more filters for feature detection. Figure 2.5 shows an input matrix is convolved with a filter. The filter moves over the entire input matrix, performing an element-wise multiplication between the input pixels  $I$  and the filter pixels  $F$ , storing the sum of the result in the corresponding element in the feature map  $R$ .

$$\begin{array}{c}
 \begin{array}{cccc}
 8 & 2 & 4 & 4 \\
 1 & 8 & 2 & 4 \\
 4 & 1 & 8 & 2 \\
 8 & 4 & 1 & 8
 \end{array} \otimes \begin{array}{cc}
 1 & 0 \\
 0 & 1
 \end{array} = \begin{array}{cccc}
 16 & 4 & 8 \\
 2 & 16 & 4 \\
 8 & 2 & 16
 \end{array}
 \end{array}$$

Image  $I$       Filter  $F$       Result  $R$

$$R[n_i, n_j] = \sum_{i=n_i}^{n_i+F_{height}} \sum_{j=n_j}^{n_j+F_{width}} I[i, j]F[i - n_i, j - n_j]$$

**Figure 2.5:** Convolution explained using a diagonal line as a filter. The filter moves over the entire input image, each element in the input image is multiplied element-wise with the corresponding element in the feature map, and the total sum is stored as the output layer.

Mathematically, convolution is implemented using cross-correlation:  
in a CNN  
as

$$r_{I,F}[n_i, n_j] = \sum_{i=n_i}^{n_i+F_{height}} \sum_{j=n_j}^{n_j+F_{width}} I[i, j] \times F[i - n_i, j - n_j] \quad (2.7)$$

where a filter image  $F$  is passed over an image  $I$  to produce a smaller image with dimensions  $(I_{width} - F_{width} + 1, I_{height} - F_{height} + 1)$ . For example, if  $n_i = 3, n_j = 2, F_{height} = 2, F_{width} = 2$  as shown in turquoise in the figure, then:

$$\begin{aligned} r_{I,F}[3, 2] &= \sum_{i=3}^5 \sum_{j=2}^4 I[i, j] F[i - 3, j - 2] \\ r_{I,F}[3, 2] &= I[3, 2]F[0, 0] + I[3, 3]F[0, 1] + I[4, 2]F[1, 0] + I[4, 3]F[1, 1] \\ r_{I,F}[3, 2] &= 1 * 1 + 8 * 0 + 4 * 0 + 1 * 1 \\ r_{I,F}[3, 2] &= 2 \end{aligned} \quad (2.8)$$

Figure 2.6 shows a convolutional layer in more detail. It consists of an input image  $I$  with the respective RGB channels. The RGB channels are convolved with each filter  $F$ , to form convolved images  $R$ . The convolved images are then passed to an activation function similar to CNNs, resulting in a ~~convolved~~ output matrix  $V$ . This entire process is referred to as a convolutional layer.

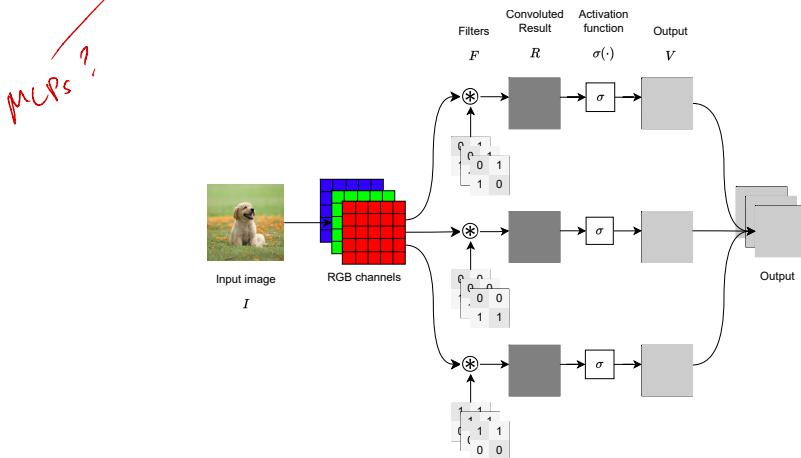


Figure 2.6

A ~~pooling layer~~ is are used to reduce the dimension of the feature maps. Similar to the convolutional layer, pooling takes a mask with fixed dimensions and aggregates the image values that fall within the mask. Max-pooling obtains the maximum value in the feature map, while average-pooling obtains the average. The fully connected layers are as described above, where the output of the convolutional layers and pooling layers are used as input to the fully connected layers.

in the previous section

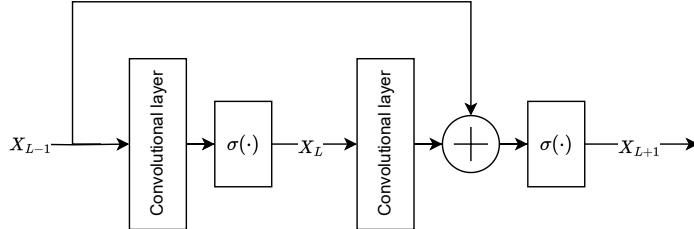
## 2.2.4. ResNet

A residual network (ResNet)

A ResNet is a convolutional neural network that aims to solve the problem of vanishing gradients which plagues deep architectures by introducing Residual Blocks. A Residual Block consists of two convolutional layers with a skip connection as shown in Figure 2.7. The skip connection allows the model to bypass some convolutional layers during training, allowing for deeper models without influencing the performance.

This allows

to be trained more successfully

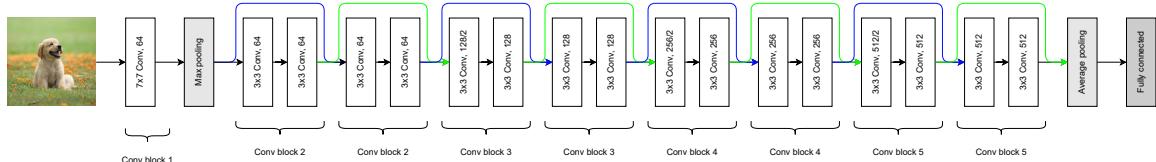


without the skip connection  
deep architectures are hard to train (the performance is worse)

**Figure 2.7:** A resnet convolutional block showing the skip connection.

ResNet

A Residual network consists of several of these Residual blocks in series. Figure 2.8 shows the general layout of a ResNet18 model (18 refers to the total number of connected layers). The model starts with a  $7 \times 7$  convolutional layer followed by a max pooling layer. It then consists of four residual convolutional blocks that are each repeated twice (for different ResNet architectures, the number of repeated blocks differs), where each convolutional block consists of two convolutional layers with the skip implemented as described above. The model ends with an average pooling layer and a fully connected layer.



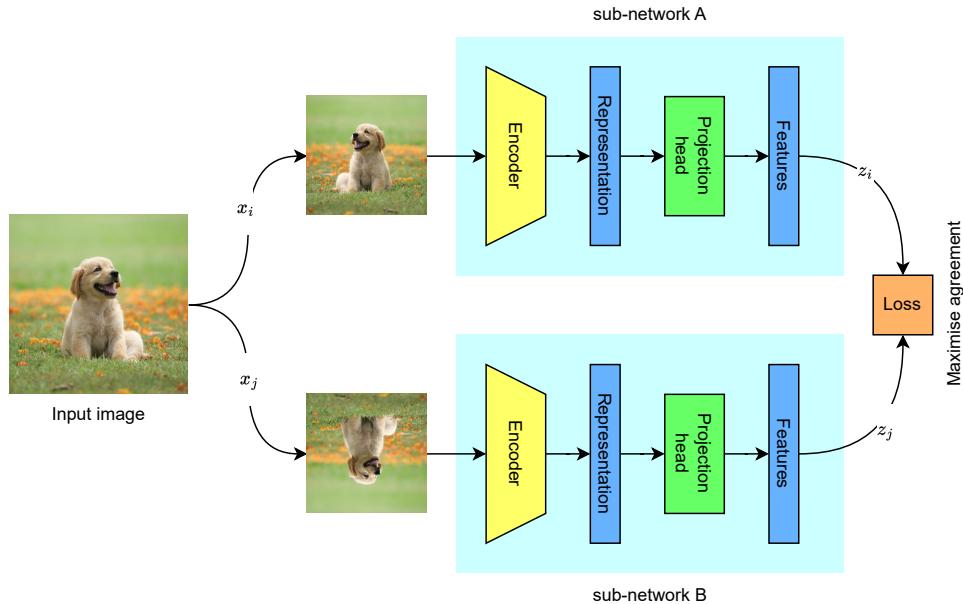
**Figure 2.8:** ResNet18 model consisting of a convolutional layer, a max pooling layer, four convolutional blocks (repeated twice), an average pooling layer and a fully connected layer.

### 2.2.5. Siamese network

A Siamese network is a neural network that consists of two sub-networks, which share the same weights. Figure 2.9 shows the layout of a basic Siamese network where the two sub-networks are shown as parallel branches. Typically, the input for this model will be a high dimensional representation, such as a mel spectrogram image, while the features generated will be a lower dimensional representation. Each sub-network consists of an encoder  $f(\cdot)$ , which returns a representation of the input, followed by a projection head  $g(\cdot)$  which returns the features  $z_i$  and  $z_j$ . The encoder model typically consists of a ResNet model to encode the input images into representations, while the projection head uses a multi-layer perceptron to generate features. These features can then be used as input to another model, for example a classifier. A contrastive loss function is used to maximise the agreement between the two input samples.

*Explain how  $x_i$  and  $x_j$  are generated.*

*explain why they should agree.*



**Figure 2.9:** A general Siamese network takes in an image  $\mathbf{x}$ , creates a copy of the original input image  $\mathbf{x}_i$  and an augmentation of the original image  $\mathbf{x}_j$ . Each sub-network consists of an encoder  $f(\cdot)$ , and a projection head  $g(\cdot)$  which generates features  $\mathbf{z}_i$  and  $\mathbf{z}_j$ . The loss function uses the features generated to maximise the agreement between the two sub-networks.

**Gradient collapse:** Gradient collapse occurs when the gradients of the loss function used in Siamese networks become very small. This is a common phenomenon in contrastive learning where loss functions are used to minimise the distance between similar images. Each of the architectures discussed below has a unique way of ensuring that the architecture does not collapse.