# Introduction to Deep Neural Networks

SSC442; Spring 2023

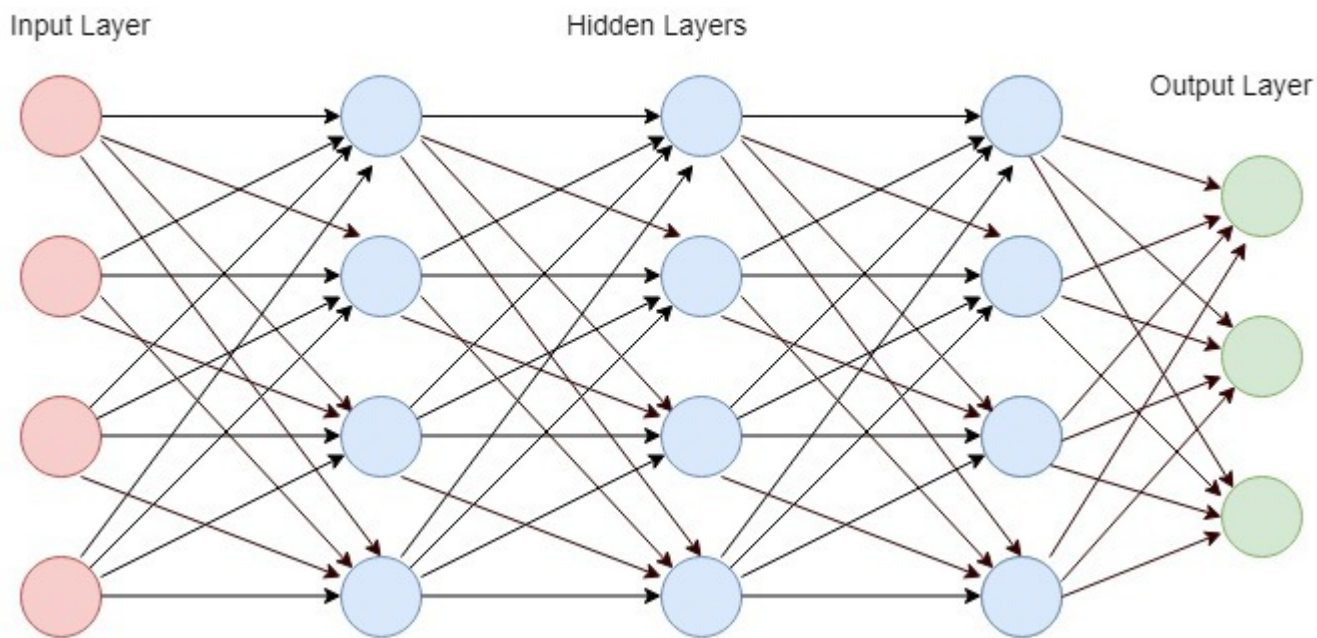Prof. Ben Bushong
Last updated March 30, 2023

# What is a Neural Network?

A neural network is a computing system inspired by the biological neural networks that constitute animal brains. It is designed to recognize patterns and learn from data.

Neural networks consist of layers of interconnected nodes called *neurons*, which process and transmit information to one another.

Neural networks have a rich history that dates back to the early 20th century. They have evolved through several significant milestones and developments.

1. **Early Concepts (1940s-1950s):** Inspired by the human brain's structure, Warren McCulloch and Walter Pitts proposed the first artificial neuron model in 1943. In 1958, Frank Rosenblatt introduced the Perceptron, a simple single-layer neural network.
2. **Backpropagation (1980s):** The backpropagation algorithm, developed by Geoffrey Hinton and others, enabled more efficient training of multi-layer neural networks. This period saw a resurgence of interest in neural networks.
3. **Deep Learning (2000s-Present):** Advances in computing power and the availability of large datasets allowed researchers to develop deeper and more complex neural networks, leading to the modern era of deep learning.

Despite facing skepticism and periods of reduced interest, neural networks have persisted and evolved to become a cornerstone of modern artificial intelligence.

Several breakthroughs and innovations have contributed to the development of neural networks over time. Some of these key developments include:

1. **Convolutional Neural Networks (CNNs):** Yann LeCun's development of the LeNet-5 architecture in the 1990s laid the foundation for modern CNNs, which are now widely used for image recognition and computer vision tasks.
2. **Recurrent Neural Networks (RNNs):** RNNs, capable of processing sequences and retaining information from previous time steps, found applications in natural language processing and time series analysis.
3. **Long Short-Term Memory (LSTM) Networks:** Hochreiter and Schmidhuber introduced LSTM networks in 1997 to address the vanishing gradient problem in RNNs, allowing for more effective learning of long-term dependencies.
4. **Transformers and Attention Mechanisms:** Vaswani et al's introduction of the Transformer architecture and attention mechanisms in 2017 revolutionized natural language processing, leading to state-of-the-art models like BERT and GPT.

These advancements have expanded the capabilities and applications of neural networks, making them a powerful tool in artificial intelligence.

# Three Primary Elements Make Up a Neural Network

1. **Input Layer:** Receives input data and passes it to the hidden layers.
2. **Hidden Layer(s):** Performs calculations and transformations on the input data.
3. **Output Layer:** Provides the final prediction or classification.

Each neuron in a layer is connected to all neurons in the previous and the next layers. These connections have weights, which determine the strength of the signal transmitted between neurons.

## Activity over the Network

Activation functions introduce non-linearity into the neural network, allowing it to learn complex patterns. Some common activation functions include:

1. **Sigmoid:** Maps input values to the range (0, 1)
2. **ReLU (Rectified Linear Unit):** Maps input values to the range [0, infinity)
3. **Tanh:** Maps input values to the range (-1, 1)
4. **Softmax:** Normalizes output values, ensuring they sum up to 1

Each neuron has an activation function that determines its output based on the weighted sum of its input signals.

# Learning in Neural Networks

A neural network learns by adjusting its weights to minimize the error between its predictions and the true labels. This process is called **training**.

1. **Forward Propagation:** Calculate predictions by passing input data through the network.
2. **Loss Calculation:** Compute the error between predictions and true labels.
3. **Backpropagation:** Update weights by propagating the error backwards through the network.
4. **Optimization:** Adjust weights using an optimization algorithm, such as gradient descent or its variants.

To ensure a neural network generalizes well to new data, it is essential to split the dataset into three parts:

1. **Training Set:** Used to train the neural network.
2. **Validation Set:** Used to tune hyperparameters and evaluate the model during training.
3. **Test Set:** Used to evaluate the model's performance after training is complete.

We avoid overfitting by monitoring the model's performance on the validation set and stopping the training process once it starts to degrade.

**MICHIGAN STATE** UNIVERSITY

Regularization techniques help prevent overfitting and improve the generalization of neural networks.

1. **L1 and L2 Regularization:** Add penalties to the loss function based on the size of the weights.
2. **Dropout:** Randomly deactivate a fraction of neurons during training, preventing co-adaptation of features.
3. **Early Stopping:** Stop training when the validation loss stops improving or starts to degrade.
4. **Batch Normalization:** Normalize input data during training, reducing the risk of vanishing or exploding gradients.

# Convolutional Neural Networks (CNNs)

CNNs are specialized neural networks designed for processing grid-like data, such as images.

1. **Convolutional Layer:** Applies a set of filters to the input, detecting local patterns.
2. **Pooling Layer:** Reduces spatial dimensions by downsampling, retaining important features while reducing computational complexity.
3. **Fully Connected Layer:** Processes the output of the convolutional and pooling layers, providing final predictions or classifications.

CNNs have shown exceptional performance in computer vision tasks, such as image classification and object detection.

We will implement a simple neural network to solve the XOR problem. Given two binary inputs, the network will predict the XOR output.

| Input 1 | Input 2 | XOR Output |
|---------|---------|------------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

We will use the Keras library in R to create our neural network.

First, we need to create the input and output data for our XOR problem. We have two input values (Input 1 and Input 2) and one output value (XOR Output).

| Input 1 | Input 2 | XOR Output |
|---------|---------|------------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

We store these values in two variables called `input_data` and `output_data`.

```
input_data <- matrix(c(0, 0, 0, 1, 1, 0, 1, 1), ncol = 2, byrow = TRUE)
output_data <- c(0, 1, 1, 0)
```

**MICHIGAN STATE** UNIVERSITY

We'll now create a simple neural network with one hidden layer and an output layer. The hidden layer has two neurons that use the ReLU activation function, which helps the network learn complex patterns.

The output layer has one neuron with a sigmoid activation function, which gives us the final prediction.

```r
model <- keras_model_sequential()

model %>%
  layer_dense(units = 2, activation = "relu", input_shape = 2) %>%
  layer_dense(units = 1, activation = "sigmoid")
```

Then, we tell the network how to learn using an optimizer, a loss function, and a metric to measure its performance.

```r
model %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = "accuracy"
)
```

**MICHIGAN STATE** UNIVERSITY

We train the neural network using the input and output data. We repeat this process 1000 times (epochs) to help the network learn better.

```
model %>% fit(
    input_data, output_data,
    epochs = 1000,
    batch_size = 4,
    validation_split = 0.25,
    verbose = 0
)
```

After training, we check how well our network learned by comparing its predictions to the true output values. We calculate the accuracy as the percentage of correct predictions.

```
predictions <- model %>% predict_classes(input_data)
accuracy <- sum(predictions == output_data) / length(output_data)
cat("Accuracy:", accuracy)
```

Our simple neural network should have an accuracy close to 100%, which means it learned to solve the XOR problem correctly.

MICHIGAN STATE UNIVERSITY

I've run the code exactly as it appears above.

The accuracy is calculated as the percentage of correct predictions made by the model. The code snippet on the previous page calculates and formats the accuracy as a percentage with two decimal places. Here, we find `accuracy` = Accuracy: 100.00%.

# Keras: A High-Level Deep Learning Library

Keras is a high-level deep learning library that allows for easy and fast prototyping of neural networks. It provides a user-friendly interface for building, training, and evaluating models with minimal code.

1. **Sequential Model:** Keras offers the `Sequential` model, which is a linear stack of layers that can be easily created and configured.
2. **Layers:** Keras provides a wide variety of predefined layers (e.g., Dense, Convolutional, LSTM) that can be added to a model with simple commands.
3. **Training and Evaluation:** Keras includes built-in functions for training, such as `fit()`, and evaluation, such as `evaluate()` and `predict()`, making it easy to train and test models.

`keras_model_sequential()` is the function used in R to create a Sequential model, which is then configured by adding layers and compiling the model.

MICHIGAN STATE UNIVERSITY

Training a neural network using Keras Sequential involves the following steps:

- **Create a Sequential Model:** Initialize a Sequential model using
  `keras_model_sequential()`.

```
model <- keras_model_sequential()
```

- **Add Layers:** Add desired layers to the model, specifying their configurations
  (e.g., number of units, activation functions).

```
model %>%
  layer_dense(units = 2, activation = "relu", input_shape = 2) %>%
  layer_dense(units = 1, activation = "sigmoid")
```

- **Compile the Model:** Specify the optimizer, loss function, and metric(s) to be used during training.

```
model %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = "accuracy"
)
```

- **Train the Model:** Use the `fit()` function to train the model on the input and output data, specifying the number of epochs, batch size, and validation split (if needed).

```
model %>% fit(
  input_data, output_data,
  epochs = 1000,
  batch_size = 4,
  validation_split = 0.25
)
```

# Neural Networks in ChatGPT and Generative Models

ChatGPT, like other generative models, relies on advanced neural networks to understand and generate human-like text. These models are trained on massive amounts of text data to learn the structure, grammar, and patterns present in human language.

1. **Deep Neural Networks:** ChatGPT uses deep neural networks with multiple layers to process and understand complex language structures.
2. **Transformers and Attention Mechanisms:** ChatGPT's architecture is based on the Transformer model, which utilizes attention mechanisms to better capture context and dependencies in the input text.

These advanced techniques enable ChatGPT to generate coherent and contextually relevant responses in a conversational setting.

Training generative models like ChatGPT involves a two-step process:

1. **Pre-training:** The model is initially trained on a large corpus of text data, such as books, articles, and web pages. This unsupervised learning helps the model learn grammar, facts, and general language patterns.
2. **Fine-tuning:** The pre-trained model is further trained on a narrower dataset, often with human-generated examples, to refine its understanding of specific tasks, like answering questions or generating text based on a given prompt.

This training process allows ChatGPT to generate human-like text and adapt to various tasks, making it a powerful tool for applications like language translation, summarization, and conversation.