

## **TASK-2**

Research study of various Deep Learning models:

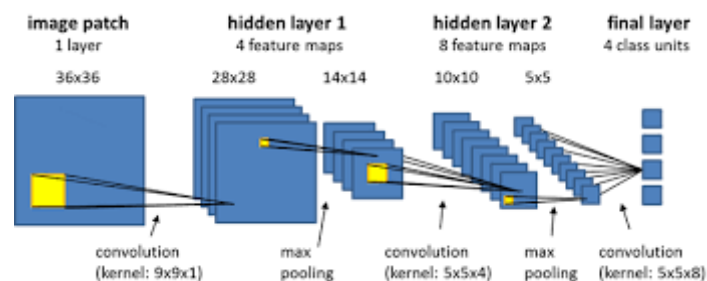
1. CNN
2. ResNet
3. Yolo (different versions)

Compare and contrast the pros and cons of each of the models and make an analysis report about the complete working and in-depth architecture. List down the features of each model which according to you are useful.

Explore other Deep Learning models.

### **Convolutional Neural Network**

Convolutional Neural Networks (CNNs) are a class of deep learning algorithms specifically designed for analyzing visual data such as images. They are widely used in computer vision tasks like image classification, object detection, and image segmentation. CNNs have a unique architecture that allows them to automatically learn and extract hierarchical representations from input images.



### **Working**

1. **Convolutional Layer:** The first layer in a CNN is the convolutional layer. It applies a set of learnable filters (also called kernels) to the input image.
2. **ReLU Activation:** ReLU introduces non-linearity into the network and helps in learning complex patterns.
3. **Pooling Layer:** Pooling layers reduce the spatial dimensions (width and height) of the feature maps while retaining the important information.

Multiple convolutional and pooling layers can be stacked together to learn more abstract and higher-level features.

4. **Fully Connected Layer:** It connects every neuron in the previous layer to the neurons in the fully connected layer.
5. **Softmax Activation:** It converts the raw scores into probability values, representing the likelihood of each class.
6. **Loss Function and Optimization:** A loss function, such as cross-entropy, is used to measure the difference between predicted and true labels. Optimization algorithms like Stochastic Gradient Descent (SGD) or Adam are commonly used to update the network parameters and minimize the loss.
7. **Backpropagation:** Backpropagation is a key technique used to update the weights of the network.

## **In-Depth Architecture**

### 1. Convolutional Layer:

The convolutional layer is the core building block of CNNs. It performs the convolution operation, which involves applying a set of learnable filters (kernels) to the input image. Each filter is a small matrix of weights that slides (convolves) across the input image, performing element-wise multiplication and summing the results to produce a feature map. The purpose of convolution is to extract local patterns and features from the image. The size and number of filters determine the depth of the feature maps generated.

### 2. ReLU Activation:

After each convolutional operation, a non-linear activation function is applied element-wise to the feature map. The Rectified Linear Unit (ReLU) activation is commonly used. ReLU sets all negative values to zero and keeps positive values unchanged. It introduces non-linearity into the network, allowing it to learn complex patterns and discriminate between different classes.

### 3. Pooling Layer:

Pooling layers are typically inserted after convolutional layers to reduce the spatial dimensions of the feature maps while retaining the important information. The most common pooling operation is max pooling, which partitions the feature map into non-overlapping rectangular regions (pooling windows) and selects the maximum value within each region. By discarding the non-maximum values, the pooling operation reduces the size of the feature maps, making them more computationally efficient. Additionally, pooling helps in creating spatial invariance, making the network more robust to translations and distortions in the input.

### 4. Convolutional and Pooling Layers (Optional):

CNNs can have multiple convolutional and pooling layers stacked together. This allows the network to learn more abstract and higher-level features by progressively combining and refining the information captured in earlier layers. Each subsequent layer learns to recognize more complex patterns and compositions of features. The deeper layers capture high-level semantic information, which is crucial for tasks like object recognition.

### 5. Fully Connected Layer:

The fully connected layer follows the convolutional and pooling layers. It is similar to the traditional neural networks where every neuron in the previous layer is connected to each neuron in the fully connected layer. The purpose of this layer is to learn the task-specific information and make predictions based on the learned features. The number of neurons in the fully connected layer corresponds to the number of classes in the classification task.

### 6. Softmax Activation:

In classification tasks, a softmax activation function is applied to the output of the fully connected layer. Softmax converts the raw scores (logits) produced by the network into probability values, representing the likelihood of each class. The sum of the probabilities for all classes is 1. The class with the highest probability is chosen as the predicted class label.

### 7. Loss Function and Optimization:

A loss function is used to measure the difference between the predicted and true labels. Cross-entropy loss is commonly used in classification tasks with softmax activation. The network aims to minimize this loss during the training process. Optimization algorithms such as Stochastic Gradient Descent (SGD) or Adam are used to update the network parameters (weights and biases) based on the gradients of the loss function. The objective is to find the optimal set of parameters that minimize the loss.

#### 8. Backpropagation:

Backpropagation is a fundamental technique used to update the weights of the network during training. It calculates the gradient of the loss function with respect to each parameter in the network and propagates it backward through the layers. The gradients are computed using the chain rule of derivatives. By iteratively adjusting the weights based on the gradients, the network learns to improve its performance and make more accurate predictions.

### **Features of CNNs:**

- Local feature extraction: CNNs are able to extract local features from input images, capturing patterns and spatial relationships.
- Hierarchical representation learning: CNNs learn hierarchical representations of data by stacking convolutional and pooling layers, enabling them to understand features at different levels of abstraction.
- Parameter sharing: CNNs share weights across different spatial locations, reducing the number of parameters and making them more efficient.
- Translation invariance: Pooling layers introduce translation invariance, allowing CNNs to handle slight variations in object position or scale.
- Generalization: CNNs demonstrate excellent generalization capabilities, able to recognize objects under different conditions and variations.

### **Pros of CNNs:**

- Superior performance in image analysis: CNNs have achieved state-of-the-art results in various computer vision tasks, outperforming traditional techniques and other machine learning algorithms.
- Automatic feature extraction: CNNs can automatically learn and extract relevant features from raw input data, eliminating the need for manual feature engineering.
- Reduced parameter space: CNNs leverage weight sharing and spatial pooling, reducing the number of parameters compared to fully connected networks.

### **Cons of CNNs:**

- Computationally intensive: Training and evaluating CNNs can be computationally demanding, requiring significant computational resources.
- Large data requirements: CNNs generally require a large amount of labeled training data to achieve good performance and avoid overfitting.
- Lack of interpretability: CNNs are often considered black-box models, making it challenging to interpret and explain their decision-making process.

## **ResNet**

### **Working of ResNet Algorithm:**

The ResNet (Residual Network) algorithm introduces skip connections, also known as residual connections, to address the vanishing gradient problem and improve the training of deep neural networks. It allows the network to learn residual mappings, making it easier to optimize deeper networks.

### **In-depth Architecture of ResNet:**

- The architecture of ResNet consists of residual blocks that contain stacked convolutional layers. Each residual block has a skip connection that bypasses one or more convolutional layers. The input to a residual block is added to the output of the block through the skip connection. This forms a shortcut path for the gradients to flow directly, allowing easier optimization of deep networks.
- The basic building block of ResNet is the residual unit. A residual unit typically contains two or three convolutional layers followed by batch normalization and ReLU activation. The output of the final convolutional layer is added element-wise to the input through the skip connection. The added input and output are then passed through another activation function to produce the final output of the residual unit.
- To construct a deeper network, multiple residual units are stacked together. The architecture can have different variations depending on the number of residual units and the size of the convolutional filters used. The common variants include ResNet-18, ResNet-34, ResNet-50, ResNet-101, and ResNet-152, where the number indicates the total number of layers in the network.
- The skip connections in ResNet allow the network to learn residual mappings, meaning it can learn to represent the difference between the input and the desired output. By using these residual connections, the network can effectively propagate gradients through the network, even in very deep architectures. This enables the training of deeper networks without encountering the vanishing gradient problem.
- ResNet architectures have revolutionized the field of deep learning, particularly in computer vision tasks. They have achieved state-of-the-art performance in various image recognition and object detection tasks, surpassing the performance of previous networks.

### **Features of ResNet:**

- Skip connections (residual connections) allow for learning residual mappings.
- Enables the training of extremely deep neural networks.
- Addresses the vanishing gradient problem.
- Improved optimization and convergence in deep networks.

### **Pros of ResNet:**

- State-of-the-art performance in various tasks, such as image recognition and object detection.
- Enables the training of much deeper networks compared to traditional architectures.
- Helps overcome the vanishing gradient problem by using skip connections.
- Residual mappings allow for better optimization and faster convergence.

### **Cons of ResNet:**

- Increased computational complexity and memory requirements due to deeper architectures.
- May suffer from overfitting if not properly regularized.
- More challenging to interpret and visualize due to the presence of skip connections.
- Training deeper networks can be more time-consuming.

## **YOLO**

### **Working of YOLO Algorithm:**

The YOLO (You Only Look Once) algorithm is an object detection algorithm that divides the input image into a grid and predicts bounding boxes and class probabilities for objects within each grid cell. It performs these predictions in a single forward pass of a convolutional neural network, making it fast and efficient for real-time object detection.

### **In-depth Architecture of YOLO:**

- The architecture of YOLO consists of a backbone network, often referred to as Darknet, which is a convolutional neural network designed for feature extraction. The backbone network typically contains multiple convolutional and pooling layers that capture and encode the hierarchical features of the input image.
- The last layers of the Darknet backbone are fully connected layers, followed by a final convolutional layer that predicts the bounding boxes and class probabilities. This final layer is responsible for dividing the input image into a grid and making predictions for each grid cell.
- Each grid cell predicts a fixed number of bounding boxes, along with the corresponding class probabilities for the objects present in that grid cell. The predicted bounding boxes consist of coordinates (x, y, width, height) relative to the grid cell and are further refined to improve accuracy.
- To improve detection accuracy, YOLO uses anchor boxes, which are predefined bounding box shapes of different aspect ratios. Each bounding box prediction is adjusted based on the anchor boxes to better match the objects in the image.
- YOLO also uses a technique called non-maximum suppression (NMS) to filter out redundant bounding box predictions. NMS suppresses overlapping bounding boxes based on their confidence scores and selects the ones with the highest scores as the final predictions.
- The architecture of YOLO allows it to detect objects directly on the full image in a single pass, without the need for region proposals or sliding windows. This design choice makes YOLO fast and efficient, suitable for real-time applications that require quick object detection.
- YOLO has gone through different versions, with each version refining the architecture and incorporating improvements in accuracy and speed. Examples of YOLO versions include

YOLOv1, YOLOv2, YOLOv3, and YOLOv4, each with its own architectural enhancements and performance optimizations.

### **Features of YOLO:**

- Fast and efficient object detection algorithm.
- Performs object detection in a single pass of a convolutional neural network.
- Divides the input image into a grid and predicts bounding boxes and class probabilities for each grid cell.
- Utilizes anchor boxes to improve the accuracy of bounding box predictions.
- Uses non-maximum suppression to filter out redundant bounding box predictions.

### **Pros of YOLO:**

- Real-time object detection: YOLO's single-pass architecture enables it to achieve real-time performance, making it suitable for applications requiring fast object detection.
- High accuracy: YOLO has shown competitive accuracy in object detection tasks, particularly for detecting objects of different sizes and aspect ratios.
- Simplicity: YOLO's unified architecture simplifies the object detection pipeline by eliminating the need for complex region proposal networks or sliding window techniques.
- End-to-end training: YOLO can be trained end-to-end, allowing for efficient optimization and faster convergence.

### **Cons of YOLO:**

- Lower localization accuracy: YOLO's grid-based approach may lead to lower precision in localizing objects, especially for small or densely packed objects.
- Difficulty in detecting small objects: Due to the grid cell size, YOLO may struggle to accurately detect small objects within an image.
- Challenging with overlapping objects: YOLO may encounter challenges when detecting and distinguishing overlapping objects that share the same grid cell.
- Limited flexibility in aspect ratios: YOLO's predefined anchor boxes may have limitations in representing objects with extreme aspect ratios or unusual shapes.