



**PERIYAR
MANIAMMAI**

INSTITUTE OF SCIENCE & TECHNOLOGY

(Deemed to be University)

Established Under Sec. 3 of UGC Act, 1956 • NAAC Accredited

think • innovate • transform

DEEP LEARNING – ASSIGNMENT 1

NAME: A.M.AKSHAYA

REGISTER NUMBER: 121011012718

COURSE NAME: DEEP LEARNING FOR DATA SCIENCE

DEPARTMENT: B.Tech C.S.E 3rd Year

LAYERS IN ARTIFICIAL NEURAL NETWORK

- 1) Input Layer: The first layer of the neural network that receives input data.

Example: In an image classification task, each pixel value of an image can be considered as an input node.

- 2) Dense Layer (Fully Connected Layer): Each neuron in this layer is connected to every neuron in the previous layer.

Example: In a feedforward neural network for image classification, a dense layer can take flattened pixel values as input.

- 3) Convolutional Layer: Contains filters that slide over input data to extract features.

Example: In a convolutional neural network (CNN) for image recognition, convolutional layers extract features like edges, textures, etc.

- 4) Pooling Layer: Reduces the spatial dimensions of the representation.

Example: Max pooling layer reduces the size of feature maps by taking the maximum value from each patch of the feature map.

- 5) Flatten Layer: Converts multi-dimensional data into a one-dimensional array.

Example: Used in CNNs to flatten the output from convolutional layers before feeding it to fully connected layers.

- 6) Dropout Layer: Randomly sets a fraction of input units to zero during training to prevent overfitting.

Example: A dropout layer with a dropout rate of 0.5 randomly sets half of

the input units to zero during training.

- 7) Batch Normalization Layer: Normalizes the activations of the previous layer at each batch.

Example: Helps in training deep networks by reducing internal covariate shift.

- 8) Activation Layer (Non-linear Activation): Introduces non-linearity into the network.

Example: ReLU (Rectified Linear Unit) activation function introduces non-linearity by outputting the input if it's positive, otherwise, it outputs zero.

- 9) Softmax Layer: Converts raw scores into probabilities.

Example: In a classification task with multiple classes, the softmax layer can convert the final layer's outputs into probability distributions over classes.

- 10) Residual Layer (Residual Block): Allows the network to learn residual functions.

Example: ResNet architecture uses residual blocks to address the vanishing gradient problem in very deep networks.

- 11) Recurrent Layer (RNN, LSTM, GRU): Processes sequences of data by maintaining internal state.

Example: LSTM (Long Short-Term Memory) layers are used in natural language processing tasks for sequential data processing.

- 12) Attention Layer: Focuses on specific parts of the input sequence.

Example: Transformer models use attention mechanisms to assign different weights to different words in a sentence.

- 13) Self-Attention Layer: An attention mechanism relating different positions of a single sequence.

Example: Used in Transformer architectures for tasks like machine translation and text generation.

- 14) Normalization Layer: Normalizes the input across the features.

Example: Layer normalization normalizes the activations of a layer across the features, rather than across the batch.

- 15) Embedding Layer: Projects categorical variables into continuous vector space.

Example: In natural language processing, words are represented as dense vectors in an embedding layer.

- 16) Concatenation Layer: Concatenates the outputs of multiple layers.

Example: In a siamese network for image similarity, the outputs of multiple convolutional layers from different branches are concatenated before making a final decision.

- 17) Addition Layer: Adds the outputs of multiple layers element-wise.

Example: Used in skip connections where the output of one layer is added to the output of another layer.

- 18) Merge Layer: Merges the outputs of multiple layers using a specific operation (e.g., addition, multiplication).

Example: Concatenating the outputs of two dense layers in a siamese network.

- 19) Gaussian Noise Layer: Adds Gaussian noise to the input.

Example: Used for regularization by adding noise to the input during training.

- 20) Global Pooling Layer: Aggregates spatial information globally.

Example: Global Average Pooling layer calculates the average value of each feature map across its entire spatial dimensions.

- 21) Local Response Normalization Layer: Normalizes the activity of neurons across adjacent layers.

Example: Used in CNN architectures like AlexNet for local contrast normalization.

- 22) Instance Normalization Layer: Normalizes the activations of each instance in a batch independently.

Example: Often used in style transfer models to normalize the activations of each instance separately.

- 23) Depthwise Separable Convolution Layer: A type of convolutional layer that factorizes a standard convolution into a depth wise convolution and a pointwise convolution.

Example: MobileNet architecture utilizes depthwise separable convolutions to reduce computational complexity.

- 24) Zero Padding Layer: Adds zero padding to the input tensor.

Example: Used in convolutional layers to preserve spatial dimensions of the input.

- 25) Spatial Dropout Layer: Randomly drops entire feature maps instead of individual elements.

Example: Used in CNNs to prevent overfitting by dropping entire feature maps during training.

- 26) Alpha Dropout Layer: Applies Alpha Dropout to the input.

Example: Used as a dropout technique with AlphaDropout activation.

- 27) Temporal Convolution Layer: Performs 1D convolution over a sequence.

Example: Used in applications like speech recognition for temporal feature extraction.

- 28) Depth Concatenation Layer: Concatenates tensors along the depth dimension.

Example: Inception modules in GoogLeNet concatenate feature maps from different convolutional paths.

- 29) Swish Activation Layer: Applies the Swish activation function element-wise.

Example: $\text{Swish}(x) = x * \text{sigmoid}(x)$.

- 30) Thresholded ReLU Layer: Applies thresholding to ReLU.

Example: Thresholded ReLU activation function sets all values below a certain threshold to zero.

- 31) Parametric ReLU Layer: ReLU with learnable parameters.

Example: PReLU (Parametric ReLU) allows the slope of the negative part of the activation to be learned.

- 32) Exponential Linear Unit (ELU) Layer: Applies ELU activation function.

Example: $\text{ELU}(x) = x$ if $x > 0$, else $\alpha * (\exp(x) - 1)$ where α is a hyperparameter.

- 33) Scaled Exponential Linear Unit (SELU) Layer: Applies SELU activation function.

Example: SELU is a variant of ELU that preserves the mean and variance of the input during training.

- 34) Scaled Dot-Product Attention Layer: Computes the dot products of the query and key vectors and scales them.

Example: Used in Transformer models for self-attention mechanism.

- 35) Multi-Head Attention Layer: Computes multiple attention mechanisms in parallel.

Example: Allows the model to focus on different parts of the input sequence simultaneously.

- 36) Positional Encoding Layer: Injects information about the relative or absolute position of the tokens in the sequence.

Example: Used in Transformer models to provide positional information to the input embeddings.

- 37) Squeeze-and-Excitation Layer: Captures interdependencies between channels.

Example: Used to recalibrate channel-wise feature responses in CNNs.

- 38) Linear Layer: Applies a linear transformation to the input data.

Example: $\text{Output} = \text{Input} * \text{Weight} + \text{Bias}$.

- 39) Gram Matrix Layer: Computes the Gram matrix of the input tensor.

Example: Used in style transfer algorithms to capture style features.

- 40) Cosine Similarity Layer: Computes the cosine similarity between two input tensors.

Example: Used in recommendation systems to measure the similarity between user preferences and item features.

- 41) Triplet Loss Layer: Computes the triplet loss for training embeddings.

Example: Used in siamese networks for learning similarity metrics.

- 42) Contrastive Loss Layer: Computes the contrastive loss for training siamese networks.

Example: Used in tasks like face verification to learn embeddings that preserve similarity/dissimilarity.

- 43) Batch-wise Contrastive Loss Layer: Computes the contrastive loss for training siamese networks on a batch-wise basis.

Example: Used in scenarios where the number of negative pairs is high, leading to computational efficiency.

- 44) Triangular Kernel Layer: Computes the triangular kernel function between two input tensors.

Example: Used in kernel methods for non-linear classification or regression tasks.

- 45) Cyclic Kernel Layer: Computes the cyclic kernel function between two input tensors.

Example: Used in cyclic-based learning algorithms.

- 46) Random Fourier Features Layer: Approximates kernel methods using random Fourier features.

Example: Used in large-scale kernel approximation techniques for efficient computation.

- 47) Zero Bias Layer: Applies a zero bias term to the input data.

Example: Used in scenarios where no bias term is required.

Identity Layer: Passes the input unchanged to the output.

Example: Used in residual connections to pass the input directly to the output.

- 48) Graph Convolutional Layer: Applies convolutional operations on graph-structured data.

Example: Used in graph neural networks for tasks like node classification, link prediction, etc.

- 49) Graph Attention Layer: Computes attention coefficients between neighboring nodes in a graph.

Example: Used in graph neural networks to prioritize information from neighboring nodes based on learned attention weights.

Types of Activation Function

Step Function:

Definition: A simple binary activation function where the output is 1 if the input is greater than or equal to a threshold, otherwise 0.

Example:

```
1 & \text{if } x \geq 0 \\
0 & \text{if } x < 0 \\
\end{cases}
```

Sigmoid Function (Logistic Function):

Definition: A smooth, S-shaped function that squashes the input values between 0 and 1.

Example: Used in binary classification problems to produce probability-like outputs.

Hyperbolic Tangent (tanh) Function:

Definition: Similar to the sigmoid function but squashes the input values between -1 and 1.

Example: Used in hidden layers of neural networks to introduce non-linearity.

Rectified Linear Unit (ReLU):

Definition: A piecewise linear function that outputs the input directly if it is positive, otherwise, it outputs zero.

Example: Widely used in deep learning due to its simplicity and effectiveness.

Leaky ReLU:

Definition: An extension of ReLU that allows a small, positive slope for negative inputs, preventing the dying ReLU problem.

Example: Helps to alleviate the vanishing gradient problem.

Parametric ReLU (PReLU):

Definition: A variant of Leaky ReLU where the slope parameter is learned during training.

Example: Allows the network to adaptively learn the slope of the activation function.

Exponential Linear Unit (ELU):

Definition: An activation function that smoothly handles negative values by allowing a small negative saturation value.

Example: Can potentially capture negative values more effectively than ReLU.

Scaled Exponential Linear Unit (SELU):

Definition: A self-normalizing activation function designed to preserve mean and variance of inputs.

Example: Used in architectures like Transformer due to its self-normalizing properties.

Softmax Function:

Definition: Converts raw scores into probabilities such that the sum of probabilities equals 1.

Example: Typically used in the output layer of a neural network for multi-class classification.

Swish Function:

Definition: A smooth, non-monotonic activation function that has been found to perform well in deep neural networks.

Example: Proposed as an alternative to ReLU with potentially better performance.

Gaussian Error Linear Unit (GELU):

Definition: A smooth approximation of the ReLU activation function based on the Gaussian cumulative distribution function.

Example: Used in transformer-based architectures for its non-linearity.

These are some of the most commonly used activation functions in neural networks, each with its own characteristics and suitability for different types of

tasks.

DIFFERENT TYPES OF OPTIMIZERS

Gradient Descent:

Definition: The basic optimization algorithm used to minimize the loss function by iteratively adjusting the model parameters in the opposite direction of the gradient.

Stochastic Gradient Descent (SGD):

Definition: A variant of gradient descent that updates the parameters using the gradient computed on a subset of the training data (mini-batch) rather than the entire dataset.

Mini-batch Gradient Descent:

Definition: An optimization algorithm that updates the model parameters using mini-batches of data, striking a balance between the efficiency of SGD and the robustness of batch gradient descent.

Momentum:

Definition: An optimization algorithm that accelerates SGD by accumulating a velocity term that carries the momentum of previous gradients.

Nesterov Accelerated Gradient (NAG):

Definition: A modification of momentum that adjusts the update direction to take into account the momentum term, resulting in faster convergence.

Adagrad (Adaptive Gradient Algorithm):

Definition: An adaptive learning rate optimization algorithm that scales the learning rate for each parameter based on the accumulated squared gradients.

RMSprop (Root Mean Square Propagation):

Definition: An adaptive learning rate optimization algorithm that divides the learning rate by an exponentially decaying average of squared gradients, preventing the learning rate from diminishing too quickly.

Adam (Adaptive Moment Estimation):

Definition: An optimization algorithm that combines the ideas of momentum and RMSprop, using both the first and second moments of the gradients to adaptively update the learning rate for each parameter.

AdaDelta:

Definition: An extension of Adagrad that dynamically adapts the learning rate over time without the need for manual tuning of the global learning rate.

AdamW:

Definition: An extension of Adam that incorporates weight decay regularization, resulting in better generalization performance.

Adamax:

Definition: A variant of Adam that replaces the second moment with the infinity norm of the gradients, making it more robust in the presence of large gradients.

Nadam (Nesterov-accelerated Adaptive Moment Estimation):

Definition: A combination of Nesterov momentum and Adam optimization, which combines the advantages of both methods.

AMSGrad:

Definition: An optimization algorithm that addresses the convergence issues of Adam by ensuring that the second moment estimate is monotonically increasing.

Adaptive Moment Estimation with Riemannian Adaptive Learning Rates (RMSpropRSL):

Definition: An optimization algorithm that adapts the learning rates on the Riemannian manifold instead of the Euclidean space, suitable for optimization problems on manifolds.

Adaptive Gradient Methods for Online Learning (AdagradRDA):

Definition: An adaptation of Adagrad that uses a per-coordinate learning rate and incorporates a regularization term to prevent divergence.

These optimizers offer different trade-offs in terms of convergence speed, memory usage, robustness to noise, and adaptability to different types of data and architectures.

