

NEURAL NETWORKS AND DEEP LEARNING

Neural Networks (NNs) are at the forefront of modern Artificial Intelligence and Deep Learning, driving advances in fields such as computer vision, natural language processing, and time-series forecasting. This document explores core neural network architectures—including Radial Basis Function (RBF) networks, Multi-Layer Perceptrons (MLPs), Long Short-Term Memory (LSTM) models, Convolutional Neural Networks (CNNs), and Recurrent Neural Networks (RNNs)—and examines their mathematical foundations, training processes, and computational complexities. By analyzing each architecture’s implementation details, advantages, and limitations, this work provides a comprehensive overview of deep learning techniques, laying the groundwork for further research and practical applications.

This document is an extension of the research and lecture notes completed at Johns Hopkins University, Whiting School of Engineering, Engineering for Professionals, Artificial Intelligence Master’s Program, Computer Science Master’s Program and Data Science Master’s Program.

Contents

1	Introduction to Neural Networks and Deep Learning	1
2	Mathematical Foundations for Neural Networks and Deep Learning	3
2.1	Mathematical Model of a Neuron	3
2.2	Activation Functions	3
2.3	Forward Propagation	3
2.4	Loss Functions	3
2.5	Backpropagation and Gradient Descent	4
2.6	Regularization Techniques	4
2.7	Summary – Mathematical Foundations for Neural Networks and Deep Learning . .	4
3	Radial Basis Function (RBF) Neural Network	5
3.1	Mathematical Formulations	5
3.2	Radial Basis Function (RBF) Neural Network without a bias term	6
3.3	Radial Basis Function (RBF) Neural Network with a bias Term	7
3.4	Radial Basis Function (RBF) Neural Network for Multi-Class Classification	9
3.5	Implementation Details	11
3.6	Advantages and Limitations	11
3.7	Analysis of the Radial Basis Function (RBF) Neural Network for Classification . . .	11
3.8	Time Complexity Analysis	11
3.9	Correctness Proof	13
3.10	Summary for Radial Basis Function (RBF) Neural Network for Classification	13
4	Multi-Layer Perceptron	15
4.1	Mathematical Formulations	15
4.2	Algorithm Explanation	16
4.3	Implementation Details	16
4.4	Advantages and Limitations	18
4.5	Analysis of the Multi-Layer Perceptron Algorithm for Classification	19
4.6	Correctness Proof	20
4.7	Summary for Multi-Layer Perceptron for Classification	21
5	Long Short-Term Memory (LSTM)	22
5.1	Mathematical Formulations	22
5.2	Algorithm Explanation	23
5.3	Implementation Details	24
5.4	Advantages and Limitations	25
5.5	Time Complexity Analysis	25
5.6	Training Complexity	25
5.7	Correctness Proof	26
5.8	Summary - LSTM	27
6	Convolutional Neural Networks	28
6.1	Architecture and Mathematical Formulation	28
6.2	Convolutional Neural Network Algorithm	30

6.3	Implementation Details	33
6.4	Advantages and Limitations	34
6.5	Analysis of the CNN Algorithms for Training and Testing	34
6.6	Correctness Proof	35
6.7	Summary for Convolutional Neural Networks	37
7	Recurrent Neural Networks (RNNs)	38
7.1	Mathematical Formulations	38
7.2	Description of the Algorithm	39
7.3	Implementation Details	41
7.4	Advantages and Limitations	42
7.5	Analysis of the RNN Algorithms for Training and Testing	43
7.6	Time Complexity Analysis	44
7.7	Correctness Proof	44
7.8	Summary of Recurrent Neural Networks	45
8	Recent Advances in Neural Networks and Deep Learning	46
8.1	PyTorch: Revolutionizing Deep Learning Frameworks	46
8.2	PyTorch Example: Simple Neural Network for Classification	46
8.3	Hugging Face: Democratizing NLP with Pre-Trained Models	48
8.4	The Intersection of PyTorch and Hugging Face	48
8.5	Summary - Recent Advances in NN and DL	49
9	Neural Networks and Deep Learning using Python Packages	50
9.1	Packages for Neural Networks and Deep Learning	50
9.2	Application of Packages to Specific Neural Network Architectures	51
9.3	Advantages and Limitations of Python Deep Learning Libraries	51
9.4	Summary - Neural Networks and Deep Learning using Python Packages	51
10	Summary	53

1 Introduction to Neural Networks and Deep Learning

Neural Networks (NNs) are computational models inspired by the interconnected structure of biological neurons. These networks learn patterns directly from data by adjusting internal parameters—weights and biases—through an iterative training process. A fundamental neural network, often referred to as a Multi-Layer Perceptron (MLP), consists of layers of interconnected neurons that progressively transform raw inputs into higher-level representations. Although the core idea remains consistent, the choice of network depth, connectivity patterns, and activation functions can significantly influence a model’s capacity to generalize and its computational requirements.

This document focuses on five key architectures that illustrate the breadth of deep learning techniques:

- **Radial Basis Function (RBF) Networks:** Utilize radially symmetric activation functions to model localized data representations, often excelling on small- to medium-scale tasks where robust local features are crucial.
- **Convolutional Neural Networks (CNNs):** Employ local receptive fields and shared weights to identify spatial or spatiotemporal patterns, originally devised for image processing but now extended to diverse structured inputs.
- **Recurrent Neural Networks (RNNs):** Incorporate hidden states that carry information across sequential data (e.g., time steps), though plain RNNs may struggle with vanishing or exploding gradients.
- **Long Short-Term Memory (LSTM):** An RNN variant designed to capture long-range dependencies via gating mechanisms, enabling more stable training over extended sequences such as in text or time-series analysis.
- **Multi-Layer Perceptrons (MLPs):** Considered foundational deep learning models with fully connected layers, suitable for a wide variety of classification and regression problems but potentially less efficient for structured or sequential inputs.

Across these architectures, certain overarching concerns apply:

- **Learning from Data:** Neural networks iteratively adjust parameters based on a loss function that quantifies discrepancies between predictions and true labels.
- **Regularization and Overfitting:** Techniques such as weight decay, dropout, and batch normalization mitigate overfitting, ensuring that networks generalize beyond training data.
- **Computational Demands:** Training modern deep networks often requires specialized hardware (e.g., GPUs) and careful hyperparameter tuning.
- **Interpretability:** Although powerful, deep neural networks can behave as “black boxes,” prompting ongoing research into interpretability strategies like saliency maps and attention mechanisms.

The surge in available data, coupled with advanced optimization methods, has propelled neural networks to the forefront of machine learning. Key applications include computer vision, language modeling, robotics, and beyond. In the following sections, we delve into the mathematical foundations underlying these architectures, illustrating how concepts such as forward propagation,

gradient-based optimization, and regularization shape deep learning's effectiveness in real-world scenarios.

2 Mathematical Foundations for Neural Networks and Deep Learning

Neural networks rely on mathematical principles that guide how inputs are processed, how errors are calculated, and how parameters are updated. These concepts form the backbone of modern deep learning architectures.

2.1 Mathematical Model of a Neuron

At their core, neurons in a neural network compute a weighted sum of inputs plus a bias term, then apply a non-linear activation function:

$$\hat{y} = f\left(\sum_{i=1}^n w_i x_i + b\right), \quad (1)$$

where x_i are input features, w_i are weights, b is a bias, and $f(\cdot)$ is an activation function (e.g., ReLU, Sigmoid, Tanh). This basic computation enables neurons to learn non-linear relationships in data.

2.2 Activation Functions

Activation functions introduce non-linearity, critical for modeling complex patterns. Common choices include:

- **Sigmoid:** $\sigma(x) = \frac{1}{1+e^{-x}}$, squashes values into $(0, 1)$, often used for binary classification.
- **ReLU:** $\max(0, x)$, accelerates convergence and mitigates vanishing gradients.
- **Tanh:** $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$, ranges in $(-1, 1)$, often helpful in certain recurrent architectures.
- **Softmax:** Normalizes a vector of logits into a probability distribution, used in multi-class classification.

2.3 Forward Propagation

Forward propagation transforms inputs through successive layers, each producing activations:

$$a^{[l+1]} = f(W^{[l]}a^{[l]} + b^{[l]}), \quad (2)$$

where $a^{[l]}$ represents the activations at layer l , and $W^{[l]}, b^{[l]}$ are layer-specific parameters. The final layer's output may be passed through a softmax for classification tasks.

2.4 Loss Functions

Loss functions quantify how well predictions match targets:

- **Mean Squared Error (MSE):** Common in regression,

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2.$$

- **Cross-Entropy Loss:** Widely used in classification,

$$\mathcal{L} = - \sum_{i=1}^N y_i \log \hat{y}_i.$$

2.5 Backpropagation and Gradient Descent

Learning occurs via backpropagation, which computes parameter gradients by applying the chain rule through each layer:

$$\frac{\partial \mathcal{L}}{\partial w_i} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_i}. \quad (3)$$

Gradient descent (or variants like Adam) updates parameters in the direction that reduces loss:

$$w \leftarrow w - \alpha \frac{\partial \mathcal{L}}{\partial w}, \quad (4)$$

where α is the learning rate, controlling the step size.

2.6 Regularization Techniques

Overfitting occurs when a model memorizes noise in the training data. Common methods to combat overfitting include:

- **L2 Regularization:** Penalizes large weights, e.g., $\mathcal{L}_{\text{reg}} = \lambda \sum_i w_i^2$.
- **Dropout:** Randomly zeroes out neuron activations, forcing more robust representations.
- **Batch Normalization:** Normalizes activations within each mini-batch, stabilizing training and improving convergence.

2.7 Summary – Mathematical Foundations for Neural Networks and Deep Learning

A strong grasp of these mathematical underpinnings empowers practitioners to design effective models. Neuron operations, activation functions, forward propagation, loss functions, backpropagation, and regularization collectively shape how neural networks learn intricate relationships. By thoughtfully applying these concepts, developers can harness the full potential of deep learning architectures, whether building simple MLPs or more sophisticated networks like CNNs, RNNs, LSTMs, or RBF-based models.

3 Radial Basis Function (RBF) Neural Network

Radial Basis Function Neural Networks (RBFNNs) are a class of supervised learning models designed for classification and regression tasks. These networks leverage radial basis functions in the hidden layer to transform input data into a high-dimensional feature space where linear separation of patterns becomes feasible. The key advantage of RBFNNs lies in their ability to model localized data representations, making them effective in handling complex, non-linear relationships and robust against noise. This property makes them particularly well-suited for multi-class classification tasks, where smooth decision boundaries are essential [2].

An RBFNN consists of three layers: an input layer that directly passes features to the next layer, a hidden layer that applies radial basis functions centered at predefined points, and an output layer that linearly combines these activations to generate predictions. The learning process involves selecting representative centers, computing their spread parameters, and solving an optimization problem to determine the optimal weights and biases in the output layer. The inclusion of a bias term in some variants enhances the model's flexibility by allowing an additional shift in the decision boundary.

This section explores various RBFNN formulations, including models with and without a bias term, and their application to multi-class classification. The mathematical formulation, training and testing procedures, computational complexity, and performance evaluation of these networks are also discussed.

3.1 Mathematical Formulations

Given a training dataset $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ and labels $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$, the RBFNN consists of three key layers:

- **Input Layer:** Passes the input features directly to the hidden layer.
- **Hidden Layer:** Applies a radial basis function centered at predefined points (centers) $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_W]^T$. The output for each center \mathbf{w}_j and input \mathbf{x}_i is:

$$H_{ij} = f\left(\frac{\|\mathbf{x}_i - \mathbf{w}_j\|}{\sigma_j}\right),$$

where σ_j is the spread of the j -th center and f is the radial basis function (e.g., Gaussian function):

$$f(r) = \exp\left(-\frac{r^2}{2}\right).$$

- **Output Layer:** Combines the hidden layer outputs linearly with weights $\hat{\mathbf{W}}$ and bias b :

$$\hat{\mathbf{y}} = \hat{\mathbf{W}}^T \mathbf{H} + b,$$

where \mathbf{H} is the vector of hidden layer activations.

The predicted class label is assigned based on:

$$\hat{y} = \arg \max_k \hat{y}[k],$$

where k indexes the output probabilities.

3.2 Radial Basis Function (RBF) Neural Network without a bias term

In this algorithm, the network omits a bias term, relying solely on weighted radial basis function responses for prediction. This section presents the training and testing algorithms for the RBF neural network. The training process involves selecting centers, computing spreads, constructing the design matrix, and solving for weights with regularization. During testing, the network predicts class labels by computing RBF activations and determining the most probable class. The absence of a bias simplifies the model while preserving expressive power through careful selection of centers and spreads. Training begins with selecting representative centers from the input data, followed by computing the spread term σ for localized basis functions to balance generalization. The final step involves solving a regularized least squares problem to determine the optimal weight matrix, ensuring robustness and improved generalization performance.

Algorithm 1 Radial Basis Function (RBF) Neural Network Training Algorithm for Classification

Require: Training dataset $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$, labels $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$, number of centers weights W , radial basis function $f(\cdot)$, regularization parameter λ .

Ensure: Trained RBF network with centers weights \mathbf{W} , weights $\hat{\mathbf{W}}$, and spreads σ .

function TRAINRBF($\mathbf{X}, \mathbf{y}, W, \phi, \lambda$)

Step 1: Select RBF Centers

Choose W centers $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_W]^T$ from \mathbf{X} using methods like k-means clustering.

Step 2: Compute Spreads

Calculate the spread σ for each center:

$$\sigma_j = \frac{1}{\sqrt{2W}} \max_{j=1, \dots, W} \|\mathbf{w}_j - \mathbf{w}_k\|, \quad \forall j, k.$$

Step 3: Construct Design Matrix

For each sample $\mathbf{x}_i \in \mathbf{X}$ and center $\mathbf{w}_j \in \mathbf{W}$, compute:

$$\mathbf{H}_{ij} = f\left(\frac{\|\mathbf{x}_i - \mathbf{w}_j\|}{\sigma_j}\right),$$

where f is the radial basis function (e.g., Gaussian).

Step 4: Solve for Weights

Add regularization to avoid overfitting:

$$\hat{\mathbf{W}} = (\mathbf{H}^T \mathbf{H} + \lambda I)^{-1} \mathbf{H}^T \mathbf{y},$$

where \mathbf{H} is the design matrix, λ is the regularization parameter, and \mathbf{y} is the label matrix.

return RBF centers \mathbf{W} , spreads σ , and weights $\hat{\mathbf{W}}$.

end function

The RBFNN for classification consists of two main phases: training and testing.

Training Phase:

1. Select W centers \mathbf{W} from the training data.
2. Compute the spreads σ for each center to control the width of the radial basis functions.
3. Construct the design matrix \mathbf{H} by computing the activations of the radial basis functions for each training sample and center.
4. Solve for the weights $\hat{\mathbf{W}}$ using regularized least squares to minimize overfitting.

Algorithm 2 Radial Basis Function (RBF) Neural Network Testing Algorithm for Classification

Require: Test sample \mathbf{x} , trained RBF network with centers \mathbf{W} , weights $\hat{\mathbf{W}}$, and spreads σ .

Ensure: Predicted class label \hat{y} .

function PREDICTRBF($\mathbf{x}, \mathbf{W}, \hat{\mathbf{W}}, \sigma$)

Step 1: Compute RBF Activations

 For each center $\mathbf{w}_j \in \mathbf{W}$, compute:

$$\mathbf{H}_j = f\left(\frac{\|\mathbf{x} - \mathbf{w}_j\|}{\sigma_j}\right).$$

Step 2: Compute Output Layer Activations

 Compute the output probabilities:

$$\hat{\mathbf{y}} = \mathbf{W}^T \mathbf{H},$$

where $\mathbf{H} = [H_1, H_2, \dots, H_W]^T$ is the vector of RBF activations.

Step 3: Assign Class Label

 Determine the predicted class:

$$\hat{y} = \arg \max_k \hat{\mathbf{y}}[k],$$

where k indexes the output probabilities.

return \hat{y} .

end function

Testing Phase:

1. Compute the radial basis function activations for the test sample using the trained centers and spreads.
2. Apply the linear combination of weights and bias to compute the class probabilities.
3. Assign the class label based on the maximum probability.

3.3 Radial Basis Function (RBF) Neural Network with a bias Term

The Radial Basis Function (RBF) Neural Network with a bias term incorporates a bias parameter in the output layer. This bias term enhances the network's flexibility, allowing for better fitting of the decision boundary by shifting activation levels appropriately. The training algorithm follows a structured approach, beginning with the selection of representative centers from the input data. The design matrix is constructed using the radial basis function activations, and an additional column of ones is appended to account for the bias term. The weight matrix, including the bias, is then obtained by solving a regularized least squares problem to mitigate overfitting.

During testing, the network predicts class labels by computing the radial basis function activations for the given input and incorporating the learned bias term. The final classification decision is made by evaluating the weighted sum of activations, including the bias, and selecting the class with the highest output probability. The inclusion of a bias term generally improves the model's adaptability and ensures more robust decision boundaries, especially in cases where the data distribution is not centered around the origin.

The RBFNN for classification consists of two main phases: training and testing.

Training Phase:

1. Select W centers \mathbf{W} from the training data.

Algorithm 3 Radial Basis Function (RBF) Neural Network Training Algorithm for Classification

Require: Training dataset $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$, labels $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$, number of centers W , radial basis function $f(\cdot)$, regularization parameter λ .

Ensure: Trained RBF network with centers \mathbf{W} , weights $\hat{\mathbf{W}}$, spreads σ , and bias b .

function TRAINRBF($\mathbf{X}, \mathbf{y}, W, f, \lambda$)

Step 1: Select RBF Centers

 Choose W centers $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_W]^T$ from \mathbf{X} .

Step 2: Compute Spreads

 Calculate the spread σ for each center:

$$\sigma_j = \frac{1}{\sqrt{2W}} \max_{j=1, \dots, W} \|\mathbf{w}_j - \mathbf{w}_k\|, \quad \forall j, k.$$

Step 3: Construct Design Matrix

 For each sample $\mathbf{x}_i \in \mathbf{X}$ and center $\mathbf{w}_j \in \mathbf{W}$, compute:

$$\mathbf{H}_{ij} = f\left(\frac{\|\mathbf{x}_i - \mathbf{w}_j\|}{\sigma_j}\right),$$

 where f is the radial basis function (e.g., Gaussian).

Step 4: Solve for Weights and Bias

 Add regularization to avoid overfitting:

$$\begin{bmatrix} \hat{\mathbf{W}} \\ b \end{bmatrix} = \left([\mathbf{H} \quad \mathbf{1}]^T [\mathbf{H} \quad \mathbf{1}] + \lambda I \right)^{-1} [\mathbf{H} \quad \mathbf{1}]^T \mathbf{y},$$

 where \mathbf{H} is the design matrix, λ is the regularization parameter, and $\mathbf{1}$ is a column vector of ones for the bias.

return RBF centers \mathbf{W} , spreads σ , weights $\hat{\mathbf{W}}$, and bias b .

end function

Algorithm 4 Radial Basis Function (RBF) Neural Network Testing Algorithm for Classification

Require: Test sample \mathbf{x} , trained RBF network with centers \mathbf{W} , weights $\hat{\mathbf{W}}$, spreads σ , and bias b .

Ensure: Predicted class label \hat{y} .

function PREDICTRBF($\mathbf{x}, \mathbf{W}, \hat{\mathbf{W}}, \sigma, b$)

Step 1: Compute RBF Activations

 For each center $\mathbf{w}_j \in \mathbf{W}$, compute:

$$\mathbf{H}_j = f\left(\frac{\|\mathbf{x} - \mathbf{w}_j\|}{\sigma_j}\right).$$

Step 2: Compute Output Layer Activations

 Compute the output probabilities:

$$\hat{\mathbf{y}} = \hat{\mathbf{W}}^T \mathbf{H} + b,$$

 where $\mathbf{H} = [H_1, H_2, \dots, H_W]^T$ is the vector of RBF activations.

Step 3: Assign Class Label

 Determine the predicted class:

$$\hat{y} = \arg \max_k \hat{\mathbf{y}}[k],$$

 where k indexes the output probabilities.

return \hat{y} .

end function

2. Compute the spreads σ for each center to control the width of the radial basis functions.
3. Construct the design matrix \mathbf{H} by computing the activations of the radial basis functions for each training sample and center.
4. Solve for the weights $\hat{\mathbf{W}}$ and bias b using regularized least squares to minimize overfitting.

Testing Phase:

1. Compute the radial basis function activations for the test sample using the trained centers and spreads.
2. Apply the linear combination of weights and bias to compute the class probabilities.
3. Assign the class label based on the maximum probability.

3.4 Radial Basis Function (RBF) Neural Network for Multi-Class Classification

In the Radial Basis Function (RBF) Neural Network multi-class setting, the output layer is designed to handle one-hot encoded labels, where each output neuron corresponds to a distinct class.

The training process involves selecting representative centers from the input data, typically using clustering methods such as k-means. Once the centers are determined, the spread parameter σ is computed to ensure appropriate generalization. The design matrix is then constructed using the radial basis function activations, and a bias term is included to enhance model expressiveness. The final weight matrix, including the bias, is obtained by solving a regularized least squares problem, ensuring robustness against overfitting.

During testing, a given input sample is transformed using the learned RBF centers and activations. The output is computed as a weighted sum of these activations, incorporating the bias term. The predicted class is determined by selecting the neuron with the highest activation, corresponding to the most probable class label. The inclusion of regularization and a bias term improves model stability and classification accuracy, making the RBF network a viable choice for multi-class problems.

The Radial Basis Function Neural Network (RBFNN) for multi-class classification operates in two main phases: training and testing. The training phase focuses on learning representative centers, computing spreads, and optimizing the weight matrix, while the testing phase applies the learned model to classify new samples.

Training Phase:

1. Select W representative centers \mathbf{W} from the training data, typically using clustering techniques such as k-means.
2. Compute the spread parameter σ for each center to control the width of the radial basis functions, ensuring a balance between local and global generalization.
3. Construct the design matrix \mathbf{H} by evaluating the radial basis function activations for each training sample with respect to the selected centers.
4. Solve for the weight matrix $\hat{\mathbf{W}}$ and bias vector \mathbf{b} using regularized least squares, incorporating a regularization parameter λ to prevent overfitting.

Testing Phase:

Algorithm 5 Radial Basis Function (RBF) Neural Network Training Algorithm for Multi-Class Classification

Require: Training dataset $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$, one-hot encoded labels $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N]^T$, number of centers W , radial basis function $f(\cdot)$, regularization parameter λ .

Ensure: Trained RBF network with centers \mathbf{W} , weights $\hat{\mathbf{W}}$, spreads σ , and bias b .

function TRAINRBF($\mathbf{X}, \mathbf{Y}, W, f, \lambda$)

Step 1: Select RBF Centers

 Choose W centers $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_W]^T$ from \mathbf{X} using methods like k-means clustering.

Step 2: Compute Spreads

 Calculate the spread σ for each center:

$$\sigma_j = \frac{1}{\sqrt{2W}} \max_{j=1, \dots, W} \|\mathbf{w}_j - \mathbf{w}_k\|, \quad \forall j, k.$$

Step 3: Construct Design Matrix

 For each sample $\mathbf{x}_i \in \mathbf{X}$ and center $\mathbf{w}_j \in \mathbf{W}$, compute:

$$\mathbf{H}_{ij} = f\left(\frac{\|\mathbf{x}_i - \mathbf{w}_j\|}{\sigma_j}\right),$$

where f is the radial basis function (e.g., Gaussian).

Step 4: Solve for Weights and Bias

 Add regularization to avoid overfitting:

$$\begin{bmatrix} \hat{\mathbf{W}} \\ \mathbf{b} \end{bmatrix} = \left(\begin{bmatrix} \mathbf{H} & \mathbf{1} \end{bmatrix}^T \begin{bmatrix} \mathbf{H} & \mathbf{1} \end{bmatrix} + \lambda I \right)^{-1} \begin{bmatrix} \mathbf{H} & \mathbf{1} \end{bmatrix}^T \mathbf{Y},$$

where \mathbf{H} is the design matrix, \mathbf{Y} contains one-hot encoded labels, λ is the regularization parameter, and $\mathbf{1}$ is a column vector of ones for the bias.

return RBF centers \mathbf{W} , spreads σ , weights $\hat{\mathbf{W}}$, and bias \mathbf{b} .

end function

Algorithm 6 Radial Basis Function (RBF) Neural Network Testing Algorithm for Multi-Class Classification

Require: Test sample \mathbf{x} , trained RBF network with centers \mathbf{W} , weights $\hat{\mathbf{W}}$, spreads σ , and bias \mathbf{b} .

Ensure: Predicted class label \hat{y} .

function PREDICTRBF($\mathbf{x}, \mathbf{W}, \hat{\mathbf{W}}, \sigma, \mathbf{b}$)

Step 1: Compute RBF Activations

 For each center $\mathbf{w}_j \in \mathbf{W}$, compute:

$$\mathbf{H}_j = f\left(\frac{\|\mathbf{x} - \mathbf{w}_j\|}{\sigma_j}\right).$$

Step 2: Compute Output Layer Activations

 Compute the output vector:

$$\hat{\mathbf{y}} = \hat{\mathbf{W}}^T \mathbf{H} + \mathbf{b},$$

where $\mathbf{H} = [H_1, H_2, \dots, H_W]^T$ is the vector of RBF activations.

Step 3: Assign Class Label

 Determine the predicted class:

$$\hat{y} = \arg \max_k \hat{y}[k],$$

where k indexes the output vector components.

return \hat{y} .

end function

1. Compute the radial basis function activations for a given test sample based on the trained centers \mathbf{W} and spreads σ .
2. Compute the output class scores by applying a linear transformation using the learned weights $\hat{\mathbf{W}}$ and bias \mathbf{b} .
3. Determine the predicted class label by selecting the class corresponding to the highest activation score.

3.5 Implementation Details

Data Preprocessing: Data preprocessing is a crucial step in machine learning to ensure consistency and improve model performance. One important aspect is to normalize or standardize input features, ensuring that they are on a consistent scale. Additionally, for multi-class classification tasks, categorical labels should be encoded as one-hot vectors to facilitate proper model training.

Model Training: During model training, k-means clustering is employed to select centers \mathbf{W} that effectively represent the data distribution. To prevent overfitting, the weight computation is regularized using a parameter λ , ensuring better generalization of the model.

Model Evaluation: Model performance is assessed using various metrics, including accuracy, precision, recall, and F1-score. Additionally, a validation set is utilized to fine-tune hyperparameters such as W and λ , ensuring optimal model performance.

3.6 Advantages and Limitations

Advantages: This approach is particularly efficient for small- to medium-sized datasets due to its localized representation. It is also robust to noise and outliers, as it leverages localized radial basis functions. Additionally, the architecture remains simple and interpretable, making it easier to analyze and understand.

Limitations: One major limitation of this approach is its computational expense for large datasets due to the design matrix computation. Additionally, its performance is highly dependent on the choice of centers and spreads, necessitating careful tuning. Furthermore, compared to deep learning models, it has limited scalability when applied to high-dimensional data.

3.7 Analysis of the Radial Basis Function (RBF) Neural Network for Classification

Radial Basis Function Neural Networks (RBFNNs) are supervised learning algorithms that employ radial basis functions as activation functions to transform the input data into a higher-dimensional space, where linear separation is feasible. This approach is widely used for multi-class classification tasks, leveraging the ability of radial basis functions to model localized patterns effectively.

3.8 Time Complexity Analysis

The time complexity of the RBFNN algorithm can be analyzed in terms of its training and testing phases:

Training Phase

1. **Selecting Centers:** Selecting W centers using k-means clustering involves:

$$O(T_k \cdot N \cdot W),$$

where T_k is the number of k-means iterations, N is the number of training samples, and W is the number of centers.

2. **Computing Spreads:** Calculating the spreads for W centers involves pairwise distance computations, resulting in:

$$O(W^2).$$

3. **Constructing Design Matrix:** Computing the design matrix \mathbf{H} requires evaluating W radial basis functions for each of the N training samples:

$$O(N \cdot W).$$

4. **Solving for Weights and Bias:** Solving the regularized least squares equation involves matrix multiplications and inversion. For N samples and W centers, the complexity is:

$$O(W^3 + N \cdot W^2).$$

Combining these, the overall training complexity is approximately:

$$O(T_k \cdot N \cdot W + W^2 + N \cdot W + W^3 + N \cdot W^2).$$

Testing Phase

1. **Computing RBF Activations:** For a single test sample, computing W radial basis function activations requires:

$$O(W).$$

2. **Output Computation:** Calculating the linear combination of weights, activations, and bias for C output classes is:

$$O(C \cdot W).$$

For M test samples, the total testing complexity is:

$$O(M \cdot W + M \cdot C \cdot W).$$

Summary - Time Complexity

The training phase has a complexity of:

$$O(T_k \cdot N \cdot W + W^3 + N \cdot W^2),$$

and the testing phase has a complexity of:

$$O(M \cdot C \cdot W).$$

These complexities highlight that the number of centers W significantly impacts both training and testing.

3.9 Correctness Proof

The correctness of the RBFNN for classification is based on its ability to approximate functions in a higher-dimensional space using localized radial basis functions.

Proof Outline

1. **Feature Transformation:** The hidden layer transforms the input data into a higher-dimensional space using radial basis functions centered at \mathbf{W} . Each function provides localized responses based on the distance from the center.
2. **Linear Combination:** The output layer combines these responses linearly, weighted by $\hat{\mathbf{W}}$ and adjusted by a bias b , ensuring the output can represent complex decision boundaries.
3. **Optimization:** The regularized least squares solution minimizes the error between predicted and actual outputs while penalizing large weights to prevent overfitting. This guarantees a solution that generalizes well to unseen data.
4. **Multi-Class Classification:** For multi-class problems, the output layer computes probabilities for each class, and the predicted class is assigned based on the maximum probability. The use of radial basis functions ensures smooth and localized decision boundaries.

Convergence Criteria

The algorithm converges when certain conditions are met. First, the number of centers W must be sufficient to capture the complexity of the data. Additionally, the regularization parameter λ should be appropriately chosen to balance bias and variance. Lastly, the optimization problem for weights and bias must be solved accurately to ensure stable convergence.

Summary - Correctness Proof

The RBFNN is correct under the assumption that radial basis functions can adequately transform the input data into a space where linear separability is achievable. The regularized optimization ensures robustness and generalization.

3.10 Summary for Radial Basis Function (RBF) Neural Network for Classification

Radial Basis Function Neural Networks (RBFNNs) provide a powerful and interpretable approach to classification by leveraging radial basis functions to transform input data into a higher-dimensional space where linear separation becomes possible. These networks are particularly effective in handling multi-class classification problems due to their ability to model complex, non-linear decision boundaries while maintaining robustness against noise.

The RBFNN training phase involves selecting representative centers using clustering techniques such as k-means, computing spread parameters to balance generalization, constructing a design matrix to capture feature transformations, and solving a regularized least squares optimization problem to determine the optimal weight matrix and bias. The testing phase applies the learned model to new data by computing radial basis function activations and assigning class labels based on the maximum activation score.

One of the key advantages of RBFNNs is their computational efficiency for small- to medium-sized datasets, owing to their localized representation. Additionally, their simple architecture and interpretable nature make them a valuable choice for classification tasks where explainability is crucial. However, RBFNNs require careful hyperparameter selection, including the number of centers and the regularization parameter, to ensure optimal performance. Furthermore, while effective for lower-dimensional data, their scalability is limited compared to deep learning models when applied to high-dimensional datasets.

Overall, RBFNNs remain a robust and versatile tool for classification, particularly in scenarios where non-linear relationships need to be captured efficiently. Their structured training approach and ability to create smooth decision boundaries make them a competitive alternative to other machine learning techniques [21].

4 Multi-Layer Perceptron

The Multi-Layer Perceptron (MLP) is a fundamental type of artificial neural network widely used for supervised learning tasks, particularly classification and regression. It consists of multiple layers of neurons, including an input layer, one or more hidden layers, and an output layer. Each neuron is connected to the next layer through weighted connections and biases, with activation functions introducing non-linearity to the model. The key advantage of MLPs is their ability to approximate complex, non-linear functions, making them suitable for diverse real-world applications such as image recognition, natural language processing, and medical diagnosis [24].

MLPs are trained using the backpropagation algorithm, which consists of two main phases: the forward pass and the backward pass. In the forward pass, input data propagates through the network, where each layer applies an activation function to transform the data before passing it to the next layer. The final layer produces class probabilities using the softmax function. The backward pass involves computing the gradient of a loss function (e.g., cross-entropy loss) and updating the network's weights using an optimization algorithm such as stochastic gradient descent (SGD) or Adam.

Despite their strong learning capabilities, MLPs require careful tuning of hyperparameters such as the number of hidden layers, neurons per layer, learning rate, and regularization techniques to prevent overfitting. Compared to simpler models like logistic regression, MLPs can learn highly non-linear decision boundaries but are computationally more expensive, especially for large datasets. The following sections explore the mathematical formulation, algorithmic details, implementation, and evaluation of MLPs in classification tasks.

4.1 Mathematical Formulations

Given a training dataset $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ with labels $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$, an MLP predicts the output $\hat{\mathbf{y}}$ for each input \mathbf{x}_i .

Forward Pass: For each layer l , the pre-activation \mathbf{z}_l and activation \mathbf{a}_l are computed as:

$$\mathbf{z}_l = \mathbf{W}_l \cdot \mathbf{a}_{l-1} + \mathbf{b}_l, \quad \mathbf{a}_l = \sigma(\mathbf{z}_l),$$

where σ is the activation function, such as ReLU or sigmoid.

For the output layer L , softmax is applied to compute class probabilities:

$$\hat{\mathbf{y}} = \text{Softmax}(\mathbf{z}_L), \quad \hat{\mathbf{y}}[k] = \frac{e^{\mathbf{z}_L[k]}}{\sum_{j=1}^C e^{\mathbf{z}_L[j]}},$$

where k indexes the output classes.

Loss Function: The cross-entropy loss for a single training sample (\mathbf{x}_i, y_i) is:

$$L(\hat{\mathbf{y}}, y_i) = - \sum_{k=1}^C y_i[k] \log(\hat{\mathbf{y}}[k]),$$

where $y_i[k]$ is a one-hot encoding of the true label.

Backward Pass (Backpropagation): The gradient of the loss with respect to weights and biases is computed using the chain rule:

$$\nabla \mathbf{W}_l = \nabla L \cdot \mathbf{a}_{l-1}^T, \quad \nabla \mathbf{b}_l = \nabla L.$$

The error is propagated backward:

$$\nabla L = \mathbf{W}_l^T \cdot \nabla L \odot \sigma'(\mathbf{z}_l),$$

where σ' is the derivative of the activation function.

4.2 Algorithm Explanation

The key components of the MLP algorithm includes an input layer, one or more hidden layers with nonlinear activation functions, and an output layer that produces the final classification scores. The network learns by adjusting weights and biases through backpropagation, optimizing a loss function such as cross-entropy using gradient-based methods like stochastic gradient descent (SGD).

The MLP training process consists of two main phases: the forward pass and the backward pass. During the forward pass, the input data is propagated through the network, where each neuron applies an activation function to a weighted sum of its inputs. The final layer computes the predicted class probabilities using the softmax function. The backward pass then calculates the gradients of the loss function with respect to the model parameters using the chain rule. These gradients are used to update the weights and biases through gradient descent, allowing the network to minimize classification error over multiple training iterations.

The testing phase follows a similar procedure, except that no weight updates are performed. Given a test sample, the MLP computes activations through the forward pass and assigns a class label based on the highest probability score in the output layer. This structured approach allows MLPs to learn complex, non-linear decision boundaries and generalize well to unseen data. The following sections present the detailed training and testing algorithms for MLP-based classification.

Training Phase:

1. Perform a forward pass to compute activations and the output $\hat{\mathbf{y}}$.
2. Compute the loss gradient ∇L using backpropagation.
3. Update the weights \mathbf{W}_l and biases \mathbf{b}_l using gradient descent:

$$\mathbf{W}_l \leftarrow \mathbf{W}_l - \eta \cdot \nabla \mathbf{W}_l, \quad \mathbf{b}_l \leftarrow \mathbf{b}_l - \eta \cdot \nabla \mathbf{b}_l.$$

Testing Phase:

1. Perform a forward pass through the trained network to compute class probabilities.
2. Assign the class label $\hat{y} = \arg \max_k \hat{\mathbf{y}}[k]$.

4.3 Implementation Details

Data Preprocessing: Proper data preprocessing is essential for ensuring efficient training and improved model performance in a Multi-Layer Perceptron (MLP). First, input features should

Algorithm 7 Multi-Layer Perceptron (MLP) Training Algorithm for Classification

Require: Training dataset $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$, labels $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$, number of epochs T , learning rate η , network architecture with weights \mathbf{W} and biases \mathbf{b} .

Ensure: Trained MLP model.

function TRAINMLP($\mathbf{X}, \mathbf{y}, T, \eta, \mathbf{W}, \mathbf{b}$)

for $t = 1$ to T **do**

▷ Epoch loop

for each training sample (\mathbf{x}_i, y_i) in (\mathbf{X}, \mathbf{y}) **do**

Forward Pass:

 Compute activations for each layer l :

$$\mathbf{z}_l = \mathbf{W}_l \cdot \mathbf{a}_{l-1} + \mathbf{b}_l, \quad \mathbf{a}_l = \sigma(\mathbf{z}_l),$$

where σ is the activation function (e.g., ReLU, sigmoid).

 Compute output layer activations:

$$\hat{\mathbf{y}} = \text{Softmax}(\mathbf{z}_L),$$

where L is the output layer.

Backward Pass:

 Compute loss gradient ∇L (e.g., cross-entropy loss):

$$\nabla L = \hat{\mathbf{y}} - y_i.$$

for each layer l from L to 1 (backpropagation) **do**

 Compute gradients for weights and biases:

$$\nabla \mathbf{W}_l = \nabla L \cdot \mathbf{a}_{l-1}^T, \quad \nabla \mathbf{b}_l = \nabla L.$$

 Update gradients for the previous layer:

$$\nabla L = \mathbf{W}_l^T \cdot \nabla L \odot \sigma'(\mathbf{z}_l),$$

where \odot is the element-wise product.

end for

Weight Update:

for each layer l **do**

 Update weights and biases:

$$\mathbf{W}_l \leftarrow \mathbf{W}_l - \eta \cdot \nabla \mathbf{W}_l, \quad \mathbf{b}_l \leftarrow \mathbf{b}_l - \eta \cdot \nabla \mathbf{b}_l.$$

end for

end for

end for

return Trained MLP model with updated \mathbf{W} and \mathbf{b} .

end function

Algorithm 8 Multi-Layer Perceptron (MLP) Testing Algorithm for Classification

Require: Test sample \mathbf{x} , trained MLP model with weights \mathbf{W} and biases \mathbf{b} .

Ensure: Predicted class label \hat{y} .

function PREDICTMLP($\mathbf{x}, \mathbf{W}, \mathbf{b}$)

Forward Pass:

 Initialize input activations $\mathbf{a}_0 \leftarrow \mathbf{x}$.

for each layer l from 1 to L **do**

 Compute pre-activation and activation:

$$\mathbf{z}_l = \mathbf{W}_l \cdot \mathbf{a}_{l-1} + \mathbf{b}_l, \quad \mathbf{a}_l = \begin{cases} \sigma(\mathbf{z}_l) & \text{if } l < L, \\ \text{Softmax}(\mathbf{z}_l) & \text{if } l = L. \end{cases}$$

end for

Prediction:

 Assign the class label:

$$\hat{y} = \arg \max_k \mathbf{a}_L[k],$$

 where k indexes the output probabilities.

return \hat{y} .

end function

be normalized or standardized to accelerate convergence during optimization. Additionally, categorical labels must be encoded as one-hot vectors to allow the network to process multi-class classification tasks effectively. Finally, the dataset should be split into training, validation, and testing subsets to enable model training, hyperparameter tuning, and unbiased performance evaluation.

Model Training: The training process of the Multi-Layer Perceptron (MLP) involves optimizing the network parameters to minimize classification error. The ReLU activation function is commonly used in hidden layers to introduce non-linearity, while the softmax function is applied in the output layer for multi-class classification. Hyperparameters such as the learning rate η , the number of epochs T , and the network architecture are fine-tuned using cross-validation to achieve optimal performance. Additionally, early stopping is employed to monitor validation performance and prevent overfitting by halting training when no further improvement is observed.

Model Evaluation: The performance of the Multi-Layer Perceptron (MLP) is assessed using various evaluation metrics, including accuracy, precision, recall, and F1-score, which provide insights into classification effectiveness. Additionally, training and validation loss curves are plotted to monitor the convergence of the model, ensuring that the network is learning effectively while avoiding overfitting or underfitting.

4.4 Advantages and Limitations

Advantages: The Multi-Layer Perceptron (MLP) is a powerful neural network architecture capable of modeling complex, non-linear relationships, making it suitable for a wide range of applications. Its flexible architecture allows customization in terms of the number of layers, neurons, and activation functions, enabling adaptation to various tasks. Additionally, MLP efficiently handles multi-class classification problems by employing the softmax function in the output layer, ensuring accurate probability distributions over multiple classes.

Limitations: The Multi-Layer Perceptron (MLP) can be computationally intensive, particularly for large datasets or deep architectures, as training requires significant computational resources. Additionally, achieving optimal performance necessitates careful tuning of hyperparameters, such as the learning rate, number of layers, and number of neurons, as well as selecting an appropriate network structure. Furthermore, MLPs are prone to overfitting, especially when dealing with limited training data, requiring the use of regularization techniques such as dropout or weight decay to enhance generalization.

4.5 Analysis of the Multi-Layer Perceptron Algorithm for Classification

The time complexity of the Multi-Layer Perceptron (MLP) algorithm for classification is influenced by the number of layers, neurons per layer, training samples, and epochs.

Training Phase

The training phase involves forward and backward passes through the network for each training sample.

1. Forward Pass:

For each layer l with n_l neurons and n_{l-1} input features:

$$\mathbf{z}_l = \mathbf{W}_l \cdot \mathbf{a}_{l-1} + \mathbf{b}_l,$$

which requires $O(n_{l-1} \cdot n_l)$ operations.

Summing over all layers:

$$\text{Forward Pass Complexity} = O\left(\sum_{l=1}^L n_{l-1} \cdot n_l\right).$$

2. Backward Pass:

Computing gradients of weights and biases requires the same complexity as the forward pass:

$$\text{Backward Pass Complexity} = O\left(\sum_{l=1}^L n_{l-1} \cdot n_l\right).$$

3. Weight Updates:

Updating weights and biases for all layers adds no additional complexity beyond the backward pass.

For N training samples and T epochs, the total training complexity is:

$$O(T \cdot N \cdot \sum_{l=1}^L n_{l-1} \cdot n_l).$$

Testing Phase

During testing, only the forward pass is performed for M test samples:

$$O(M \cdot \sum_{l=1}^L n_{l-1} \cdot n_l).$$

Summary - Time Complexity

The training complexity of the MLP algorithm is $O(T \cdot N \cdot \sum_{l=1}^L n_{l-1} \cdot n_l)$, while the testing complexity is $O(M \cdot \sum_{l=1}^L n_{l-1} \cdot n_l)$. The computation scales linearly with the number of training samples, test samples, and epochs.

4.6 Correctness Proof

The correctness of the Multi-Layer Perceptron algorithm for classification is based on the universal approximation theorem, which guarantees that MLPs with sufficient neurons and layers can approximate any continuous function to arbitrary precision.

Proof Outline

1. Representation Power:

The forward pass computes transformations through weights and activation functions, enabling the MLP to learn complex, non-linear decision boundaries.

2. Loss Minimization:

The backward pass adjusts weights and biases to minimize the cross-entropy loss function:

$$L(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{i=1}^N \sum_{k=1}^C y_i[k] \log(\hat{\mathbf{y}}[k]),$$

ensuring that the predicted probabilities $\hat{\mathbf{y}}$ align with the true labels \mathbf{y} .

3. Gradient Descent Convergence:

The weights are updated using the gradient of the loss function:

$$\mathbf{W}_l \leftarrow \mathbf{W}_l - \eta \cdot \nabla \mathbf{W}_l.$$

This iterative update guarantees convergence to a local minimum of the loss function, assuming a sufficiently small learning rate η .

4. Classification Decision:

The output layer applies the softmax function, which maps logits to probabilities, ensuring that the predicted class corresponds to the maximum probability.

Generalization Guarantees:

The generalization performance depends on: - Regularization techniques (e.g., dropout, weight decay). - Proper selection of network architecture (e.g., depth, width).

Summary - Correctness Proof

The MLP algorithm is correct under the assumption that the network has sufficient capacity and that the optimization converges. The universal approximation theorem and backpropagation ensure that the network can learn to approximate the underlying data distribution and classify correctly.

4.7 Summary for Multi-Layer Perceptron for Classification

The Multi-Layer Perceptron (MLP) is a powerful and versatile neural network architecture capable of learning complex, non-linear relationships in classification tasks. It consists of multiple layers of neurons, with non-linear activation functions allowing the network to map inputs to class labels effectively. The training process involves forward propagation of input data, backpropagation of error signals, and weight updates using optimization techniques like stochastic gradient descent.

The computational complexity of MLPs depends on the number of layers, neurons per layer, and training samples. The training phase, which involves iterative weight updates, has a complexity of:

$$O(T \cdot N \cdot \sum_{l=1}^L n_{l-1} \cdot n_l),$$

where T is the number of epochs, N is the number of training samples, and L is the number of layers. The testing phase, which involves only forward propagation, has a complexity of:

$$O(M \cdot \sum_{l=1}^L n_{l-1} \cdot n_l),$$

where M is the number of test samples.

MLPs offer several advantages, including their ability to model non-linear decision boundaries, flexible architecture customization, and efficient handling of multi-class classification through the softmax function. However, they also come with challenges, such as high computational costs for deep architectures, sensitivity to hyperparameter selection, and susceptibility to overfitting if not properly regularized.

The correctness of MLPs is supported by the universal approximation theorem, which guarantees that an MLP with at least one hidden layer can approximate any continuous function given sufficient neurons. This theoretical foundation, combined with practical optimizations such as dropout, weight decay, and batch normalization, makes MLPs a valuable tool for modern classification tasks.

Overall, MLPs remain a cornerstone of deep learning and artificial intelligence, serving as the foundation for more advanced architectures such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs) [24].

5 Long Short-Term Memory (LSTM)

Long Short-Term Memory (LSTM) networks are a specialized type of recurrent neural network (RNN) designed to overcome the limitations of traditional RNNs, particularly the vanishing gradient problem [14]. LSTMs are equipped with memory cells and gating mechanisms that allow them to retain important information over long sequences, making them highly effective for sequential data tasks. Their applications span a wide range of domains, including natural language processing (NLP), speech recognition, and machine translation.

Unlike conventional RNNs, which struggle to maintain long-term dependencies, LSTMs selectively store and retrieve information using forget, input, and output gates. These components regulate the flow of information through the network, enabling efficient learning from sequential data. In this section, we delve into the mathematical foundations of LSTMs, their algorithmic structure, and their computational efficiency in training and testing phases.

5.1 Mathematical Formulations

An LSTM unit at time step t is defined by the following equations:

$$\mathbf{f}_t = \sigma(\mathbf{W}_f[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f) \quad (\text{Forget gate}) \quad (5)$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_i[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i) \quad (\text{Input gate}) \quad (6)$$

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_c[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_c) \quad (\text{Candidate cell state}) \quad (7)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t \quad (\text{Cell state update}) \quad (8)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_o[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o) \quad (\text{Output gate}) \quad (9)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (\text{Hidden state update}) \quad (10)$$

where:

- σ is the sigmoid activation function.
- \tanh is the hyperbolic tangent activation function.
- $\mathbf{W}_f, \mathbf{W}_i, \mathbf{W}_c, \mathbf{W}_o$ are the weight matrices.
- $\mathbf{b}_f, \mathbf{b}_i, \mathbf{b}_c, \mathbf{b}_o$ are the bias terms.
- \mathbf{h}_t is the hidden state at time step t .
- \mathbf{c}_t is the cell state at time step t .
- \odot represents element-wise multiplication.

For text classification, the final hidden state of the LSTM is passed through a fully connected layer followed by a softmax function:

$$\hat{\mathbf{y}} = \text{Softmax}(\mathbf{W}_{\text{out}}\mathbf{h}_T + \mathbf{b}_{\text{out}}) \quad (11)$$

where $\hat{\mathbf{y}}$ represents the predicted probability distribution over classes.

5.2 Algorithm Explanation

The LSTM algorithm is designed to process sequential data by maintaining long-term dependencies using memory cells. Unlike traditional RNNs, LSTMs utilize gating mechanisms (forget, input, and output gates) to regulate information flow [11]. This enables them to retain relevant information over long sequences, making them effective for tasks such as text classification, where context matters.

The LSTM training algorithm follows these steps:

1. Initialize LSTM parameters and hyperparameters.
2. Process input sequences through LSTM layers.
3. Compute the output and loss using cross-entropy for classification.
4. Perform backpropagation through time (BPTT) and update parameters using gradient descent.
5. Repeat until convergence.

For testing, the trained model is used to predict class labels for new input sequences by feeding them through the LSTM layers and selecting the most probable class.

Algorithm 9 LSTM Training Algorithm for Text Classification

Require: Training data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$, labels $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$, number of LSTM layers L , hidden units per layer H , learning rate η , number of epochs N_{epochs} , batch size B .

Ensure: Trained LSTM model.

function TRAINLSTM($\mathbf{X}, \mathbf{y}, L, H, \eta, N_{\text{epochs}}, B$)

Initialize parameters $\mathbf{W}_f^{(l)}, \mathbf{W}_i^{(l)}, \mathbf{W}_c^{(l)}, \mathbf{W}_o^{(l)}, \mathbf{b}_f^{(l)}, \mathbf{b}_i^{(l)}, \mathbf{b}_c^{(l)}, \mathbf{b}_o^{(l)}$ for all $l = 1, 2, \dots, L$

Initialize LSTM hidden and cell states $\mathbf{h}_t^{(l)}, \mathbf{c}_t^{(l)}$ to zeros

for epoch = 1 to N_{epochs} **do**

for each batch $(\mathbf{X}_b, \mathbf{y}_b) \subset (\mathbf{X}, \mathbf{y})$ of size B **do**

$\mathbf{h}_0^{(0)} \leftarrow \mathbf{0}$

 ▷ Initialize hidden state of first layer

$\mathbf{c}_0^{(0)} \leftarrow \mathbf{0}$

 ▷ Initialize cell state of first layer

for $t = 1$ to T **do**

 ▷ Process each time step

$\mathbf{h}_t^{(0)} \leftarrow \mathbf{x}_t$

 ▷ Input to first LSTM layer

for $l = 1$ to L **do**

 ▷ Process through LSTM layers

$\mathbf{f}_t^{(l)} \leftarrow \sigma(\mathbf{W}_f^{(l)}[\mathbf{h}_{t-1}^{(l)}, \mathbf{h}_t^{(l-1)}] + \mathbf{b}_f^{(l)})$

 ▷ Forget gate

$\mathbf{i}_t^{(l)} \leftarrow \sigma(\mathbf{W}_i^{(l)}[\mathbf{h}_{t-1}^{(l)}, \mathbf{h}_t^{(l-1)}] + \mathbf{b}_i^{(l)})$

 ▷ Input gate

$\tilde{\mathbf{c}}_t^{(l)} \leftarrow \tanh(\mathbf{W}_c^{(l)}[\mathbf{h}_{t-1}^{(l)}, \mathbf{h}_t^{(l-1)}] + \mathbf{b}_c^{(l)})$

 ▷ Candidate cell state

$\mathbf{c}_t^{(l)} \leftarrow \mathbf{f}_t^{(l)} \odot \mathbf{c}_{t-1}^{(l)} + \mathbf{i}_t^{(l)} \odot \tilde{\mathbf{c}}_t^{(l)}$

 ▷ Cell state update

$\mathbf{o}_t^{(l)} \leftarrow \sigma(\mathbf{W}_o^{(l)}[\mathbf{h}_{t-1}^{(l)}, \mathbf{h}_t^{(l-1)}] + \mathbf{b}_o^{(l)})$

 ▷ Output gate

$\mathbf{h}_t^{(l)} \leftarrow \mathbf{o}_t^{(l)} \odot \tanh(\mathbf{c}_t^{(l)})$

 ▷ Hidden state update

end for

end for

$\hat{\mathbf{y}}_b \leftarrow \text{Softmax}(\text{FullyConnected}(\mathbf{h}_T^{(L)}))$

 ▷ Output layer

 Compute loss $\mathcal{L}(\hat{\mathbf{y}}_b, \mathbf{y}_b)$

 Backpropagate the loss and update parameters using gradient descent with learning rate η

end for

end for **return** Trained LSTM model

end function

Algorithm 10 LSTM Testing Algorithm for Text Classification

Require: Trained LSTM model, test data $\mathbf{X}_{\text{test}} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M]^T$, corresponding true labels $\mathbf{y}_{\text{test}} = [y_1, y_2, \dots, y_M]^T$

Ensure: Predicted labels $\hat{\mathbf{y}}_{\text{test}}$ and accuracy score

function TESTLSTM(Trained Model, $\mathbf{X}_{\text{test}}, \mathbf{y}_{\text{test}}$)

 Initialize accuracy counter correct $\leftarrow 0$

for each sequence $\mathbf{x}_m \in \mathbf{X}_{\text{test}}$ **do**

$\mathbf{h}_0^{(0)} \leftarrow \mathbf{0}$ ▷ Initialize hidden state

$\mathbf{c}_0^{(0)} \leftarrow \mathbf{0}$ ▷ Initialize cell state

for $t = 1$ to T **do**

$\mathbf{h}_t^{(0)} \leftarrow \mathbf{x}_t$ ▷ Input to first LSTM layer

for $l = 1$ to L **do**

 Compute LSTM hidden state and cell state updates as in training

end for

end for

$\hat{y}_m \leftarrow \arg \max \text{Softmax}(\text{FullyConnected}(\mathbf{h}_T^{(L)}))$ ▷ Predict class label

if $\hat{y}_m = y_m$ **then**

 correct \leftarrow correct + 1

end if

end for

 Compute accuracy $\text{Acc} = \frac{\text{correct}}{M}$ **return** Predicted labels $\hat{\mathbf{y}}_{\text{test}}$ and accuracy Acc

end function

5.3 Implementation Details

Data Preprocessing

Effective data preprocessing is crucial for optimizing the performance of an LSTM model in text classification tasks. The first step involves converting text data into numerical representations, such as word embeddings, to enable the model to process linguistic information in a structured format. Next, text sequences are normalized and tokenized, ensuring that variations in capitalization, punctuation, and special characters do not affect the model's performance. Additionally, sequences are padded to maintain a uniform input length, allowing the LSTM to handle varying text sizes efficiently while preserving temporal dependencies across sequences.

Model Processing

The LSTM model undergoes several critical processing steps to ensure effective training and performance evaluation. First, the network is initialized with the specified number of layers and hidden units, which define the capacity of the model to capture sequential dependencies. During training, the model optimizes its parameters using stochastic gradient descent (SGD) or the Adam optimizer, both of which help in adjusting weights based on the computed loss function. Finally, the model's performance is validated using a separate test dataset, ensuring that it generalizes well to unseen data and does not overfit the training set.

Evaluation

The evaluation of an LSTM model involves assessing its performance using various quantitative metrics to ensure its effectiveness in text classification tasks. Key performance indicators include

accuracy, precision, recall, and F1-score, which provide insights into the model’s predictive capabilities and balance between false positives and false negatives. Additionally, a confusion matrix is analyzed to identify error patterns, offering a deeper understanding of misclassifications and potential areas for improvement in the model’s learning process.

5.4 Advantages and Limitations

LSTM networks offer several advantages that make them well-suited for sequential data processing. They effectively capture long-range dependencies in text sequences, allowing them to maintain context over extended input sequences. Additionally, LSTMs are highly effective in processing sequential data due to their ability to retain relevant information while mitigating the vanishing gradient problem. Furthermore, they can be stacked into multiple layers, enhancing their expressiveness and enabling them to learn hierarchical representations of complex patterns.

Despite their advantages, LSTM networks have certain limitations. They are computationally expensive due to their recurrent structure, requiring significant processing power and memory for training on long sequences. Additionally, LSTMs are prone to overfitting, particularly when trained on small datasets, necessitating the use of regularization techniques such as dropout. Furthermore, compared to convolutional neural networks (CNNs), LSTMs exhibit slower training times due to the sequential nature of their computations, which limits their parallelization efficiency.

5.5 Time Complexity Analysis

The time complexity of LSTM training and testing depends on multiple factors, including the number of input sequences, the sequence length, the number of LSTM layers, and the number of hidden units per layer.

5.6 Training Complexity

During training, for each time step t in an input sequence, each LSTM layer performs matrix multiplications, additions, and activation functions.

Given an LSTM model, the key parameters influencing its complexity and performance include N , the number of training samples, T , the sequence length, H , the number of hidden units per layer, and L , the total number of LSTM layers. These parameters collectively determine the computational requirements and the model’s ability to capture long-term dependencies in sequential data.

Each LSTM unit at time step t performs two primary computational operations. The matrix-vector multiplication, which is responsible for updating the gates and cell states, has a complexity of $O(H^2)$. Additionally, element-wise operations, such as activation functions and gating mechanisms, contribute a complexity of $O(H)$. Together, these operations determine the overall computational cost of processing sequential data within an LSTM network.

Since each sequence has length T , and we perform these computations across L layers, the per-sequence complexity is:

$$O(LTH^2). \tag{12}$$

Considering N training samples and E training epochs, the total training complexity is:

$$O(ENLTH^2). \quad (13)$$

This shows that increasing the number of layers (L), sequence length (T), or hidden units (H) significantly impacts training time.

Testing Complexity

During inference, we only perform forward passes through the LSTM without backpropagation. Each input sequence requires:

$$O(LTH^2). \quad (14)$$

For M test samples, the total testing complexity is:

$$O(MLTH^2). \quad (15)$$

Since $M \ll N$ (typically), testing is significantly faster than training.

Summary of Time Complexity

The time complexity of the LSTM algorithm varies between training and testing phases. During **training**, the complexity is $O(ENLTH^2)$, where E is the number of epochs, N is the number of training samples, L is the number of LSTM layers, T is the sequence length, and H is the number of hidden units. This complexity is higher due to the computational overhead of backpropagation through time. In contrast, during **testing**, the complexity is reduced to $O(MLTH^2)$, where M represents the number of test samples. Since testing does not involve gradient updates, it is significantly faster than training.

5.7 Correctness Proof

The correctness of the LSTM algorithm can be established by ensuring that the computations performed by the LSTM gates correctly propagate information through time.

Correctness of LSTM Computation

The LSTM cell updates its hidden and cell states using:

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t, \quad (16)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t). \quad (17)$$

Each gate: The forget gate ensures retention of past information by determining which parts of the previous cell state should be carried forward. The input gate allows updates based on new input, deciding which information should be added to the cell state. Finally, the output gate controls how much of the memory should be exposed as the hidden state, influencing the information passed to the next time step.

Correctness of Backpropagation

Gradient updates for each weight matrix in LSTM are computed using:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}}. \quad (18)$$

Since all components are differentiable and gradients are computed recursively using backpropagation through time (BPTT) graves2012supervised, the correctness of LSTM training follows.

Correctness of Classification

For text classification, the output $\hat{\mathbf{y}}$ is computed using:

$$\hat{\mathbf{y}} = \text{Softmax}(\mathbf{W}_{\text{out}} \mathbf{h}_T + \mathbf{b}_{\text{out}}). \quad (19)$$

Since the softmax function ensures that probabilities sum to 1, the predicted label is:

$$y^* = \arg \max \hat{\mathbf{y}}. \quad (20)$$

Given a sufficiently trained model, this classification follows from standard supervised learning correctness properties.

Summary of Correctness Proof

LSTM maintains long-term dependencies through gated memory updates, allowing it to effectively capture sequential patterns in data. Additionally, gradient computations ensure correct parameter updates via backpropagation, enabling the model to learn meaningful representations. Furthermore, the softmax function guarantees valid probability distributions for classification, ensuring that the model produces well-calibrated class probabilities.

5.8 Summary - LSTM

LSTM networks provide a powerful approach to processing sequential data, overcoming the vanishing gradient issue that limits traditional RNNs. By utilizing a gating mechanism, they maintain long-term dependencies and achieve high performance in tasks such as text classification, speech processing, and time-series prediction.

This section explored the mathematical formulation of LSTMs, the structure of their training and testing algorithms, and their computational efficiency. The analysis of time complexity demonstrated that while LSTMs are computationally expensive due to their recurrent nature, they remain a valuable tool for tasks requiring sequential pattern learning. Furthermore, the correctness proof established the validity of LSTM computations, ensuring reliable classification performance. Despite their strengths, LSTMs require careful hyperparameter tuning and efficient training strategies to mitigate issues such as overfitting and long training times.

Ultimately, LSTMs serve as a fundamental component in deep learning architectures for sequential data modeling, providing a robust solution for applications that demand contextual understanding over extended input sequences.

6 Convolutional Neural Networks

Convolutional Neural Networks (CNNs) have revolutionized the field of deep learning by introducing hierarchical feature extraction through convolutional layers [18]. Unlike fully connected networks, CNNs preserve spatial relationships within data by leveraging local connectivity and weight sharing, making them particularly effective for structured inputs such as images, speech, and text [17].

A CNN is composed of multiple layers, including convolutional layers, activation functions, pooling layers, and fully connected layers, all working together to automatically extract relevant features from input data [10]. The ability to learn spatial hierarchies enables CNNs to outperform traditional machine learning approaches in computer vision tasks such as image classification, object detection, and segmentation. Additionally, CNNs have been successfully applied to natural language processing (NLP), where they capture local dependencies in textual data, making them useful for tasks such as sentiment analysis and document classification [15].

Training a CNN involves forward propagation, where input data is transformed through layers of learned filters, and backpropagation, where errors are propagated backward to adjust the filter weights [23]. By optimizing a loss function using gradient-based methods such as Stochastic Gradient Descent (SGD) or Adam, CNNs can learn intricate data representations with high accuracy [16]. Despite their impressive performance, CNNs require large amounts of labeled training data and computational resources, often necessitating the use of GPUs for efficient processing [17].

This section provides a comprehensive analysis of CNNs, including their mathematical formulation, training process, algorithmic implementation, computational complexity, and correctness proof.

6.1 Architecture and Mathematical Formulation

A CNN consists of several key components:

Convolutional Layer

The convolutional layer applies a set of learnable filters to extract local patterns such as edges, textures, and shapes from the input [18]. The convolution operation preserves spatial relationships while reducing the number of parameters compared to fully connected networks. Mathematically, a convolution operation is defined as:

$$(f * w)(x, y) = \sum_{i=0}^{k-1} \sum_{j=0}^{k-1} w(i, j) \cdot x(i + p, j + q), \quad (21)$$

where $x(i, j)$ is the input, $w(i, j)$ is the filter of size $k \times k$, and (p, q) denotes the spatial location of the filter. This operation captures spatial patterns such as edges and textures.

Activation Function

Non-linearity is introduced using activation functions, typically the Rectified Linear Unit (ReLU). Non-linearities, such as the ReLU, introduce essential complexity, allowing CNNs to learn intricate hierarchical features:

$$\text{ReLU}(z) = \max(0, z). \quad (22)$$

ReLU accelerates training convergence and mitigates the vanishing gradient problem [9].

Pooling Layer

Pooling layers downsample feature maps to reduce computational complexity while preserving important information. These layers perform downsampling operations, such as max pooling or average pooling, to reduce spatial dimensions while retaining important features. This reduces computational complexity and enhances translation invariance. Max pooling is defined as:

$$y(i, j) = \max_{(p, q) \in R} x(i + p, j + q), \quad (23)$$

where R represents the receptive field of the pooling operation. This enhances translational invariance [17].

Fully Connected Layer and Output

The extracted hierarchical features are fed into fully connected layers for classification. Typically found in the final stages, these layers aggregate extracted features to make final predictions. They function similarly to those in MLPs but operate on feature representations learned by the convolutional layers:

$$y = \sigma(Wx + b), \quad (24)$$

where W and b are learnable parameters, and σ is typically the softmax function for multi-class classification.

Training and Backpropagation

CNNs are trained using stochastic gradient descent (SGD) and backpropagation, where gradients of convolutional layers are computed using the chain rule:

$$\frac{\partial L}{\partial w} = \sum_{i, j} \frac{\partial L}{\partial y(i, j)} \cdot \frac{\partial y(i, j)}{\partial w}. \quad (25)$$

Weight updates follow:

$$w \leftarrow w - \eta \frac{\partial L}{\partial w}, \quad (26)$$

where η is the learning rate [10].

6.2 Convolutional Neural Network Algorithm

The following algorithm details the essential steps in training a CNN, including forward propagation through the convolutional, activation, and pooling layers, loss computation, backpropagation for gradient updates, and optimization of network parameters.

In this formulation, the CNN processes input text representations (e.g., word embeddings) using convolutional filters to identify local patterns, followed by non-linear transformations and downsampling operations. The extracted features are then passed to a fully connected layer for classification. Below, we present the training and testing algorithms for a CNN designed for text classification.

Algorithm 11 CNN Training Algorithm for Text Classification

Require: Training dataset $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$, labels $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$, number of epochs T , learning rate η , filter weights \mathbf{F} , network parameters \mathbf{W}, \mathbf{b}

Ensure: Trained CNN model with optimized parameters

```

1: function TRAINCNN( $\mathbf{X}, \mathbf{y}, T, \eta, \mathbf{F}, \mathbf{W}, \mathbf{b}$ )
2:   for  $t = 1$  to  $T$  do ▷ Epoch loop
3:     for each training sample  $(\mathbf{x}_i, y_i)$  in  $(\mathbf{X}, \mathbf{y})$  do
4:       Step 1: Forward Pass
5:        $\mathbf{C} \leftarrow \text{CONVOLUTION}(\mathbf{x}_i, \mathbf{F})$ 
6:        $\mathbf{A} \leftarrow \text{RELU}(\mathbf{C})$ 
7:        $\mathbf{P} \leftarrow \text{MAXPOOLING}(\mathbf{A})$ 
8:        $\hat{\mathbf{y}} \leftarrow \text{FULLYCONNECTED}(\mathbf{P}, \mathbf{W}, \mathbf{b})$ 
9:       Step 2: Compute Loss
10:       $L \leftarrow \text{CROSSENTROPYLOSS}(\hat{\mathbf{y}}, y_i)$ 
11:      Step 3: Backward Pass (Gradient Calculation)
12:       $\nabla \mathbf{W}, \nabla \mathbf{b} \leftarrow \text{COMPUTEGRADIENTS}(L, \mathbf{W}, \mathbf{b})$ 
13:       $\nabla \mathbf{F} \leftarrow \text{COMPUTEFILTERGRADIENTS}(L, \mathbf{F})$ 
14:      Step 4: Parameter Update
15:       $\mathbf{W} \leftarrow \mathbf{W} - \eta \cdot \nabla \mathbf{W}$ 
16:       $\mathbf{b} \leftarrow \mathbf{b} - \eta \cdot \nabla \mathbf{b}$ 
17:       $\mathbf{F} \leftarrow \mathbf{F} - \eta \cdot \nabla \mathbf{F}$ 
18:    end for
19:  end for
20:  return Trained CNN model with updated  $\mathbf{F}, \mathbf{W}, \mathbf{b}$ 
21: end function
```

Algorithm 12 CNN CROSSENTROPYLOSS Function for Training

```

1: function CROSSENTROPYLOSS( $\hat{\mathbf{y}}, y_i$ )
2:   Compute loss:  $L = - \sum_{k=1}^C y_i[k] \log(\hat{\mathbf{y}}[k])$ 
3:   return  $L$ 
4: end function
```

Algorithm 13 CNN COMPUTEGRADIENTS Function for Training

```

1: function COMPUTEGRADIENTS( $L, \mathbf{W}, \mathbf{b}$ )
2:   Compute weight gradients:  $\nabla \mathbf{W} = \frac{\partial L}{\partial \mathbf{W}}$ 
3:   Compute bias gradients:  $\nabla \mathbf{b} = \frac{\partial L}{\partial \mathbf{b}}$ 
4:   return  $\nabla \mathbf{W}, \nabla \mathbf{b}$ 
5: end function
```

Algorithm 14 CNN COMPUTEFILTERGRADIENTS Function for Training

```
1: function COMPUTEFILTERGRADIENTS( $L, \mathbf{F}$ )
2:   Compute filter gradients:  $\nabla \mathbf{F} = \frac{\partial L}{\partial \mathbf{F}}$ 
3:   return  $\nabla \mathbf{F}$ 
4: end function
```

Algorithm 15 CNN CONVOLUTION Function for Training and Testing

```
1: function CONVOLUTION( $\mathbf{x}_i, \mathbf{F}$ )
2:   Input: Input feature map  $\mathbf{x}_i \in \mathbb{R}^{H \times W \times C}$ , filter  $\mathbf{F} \in \mathbb{R}^{K \times K \times C}$ 
3:   Output: Convolved feature map  $\mathbf{C} \in \mathbb{R}^{H' \times W' \times D}$ 
4:   Initialize  $\mathbf{C}$  as an empty tensor of size  $H' \times W' \times D$ 
5:   for each filter  $d \in \{1, \dots, D\}$  do
6:     for each spatial position  $(i, j)$  in  $\mathbf{x}_i$  do
7:       Compute convolution operation:
```

$$C_{i,j,d} = \sum_{m=0}^{K-1} \sum_{n=0}^{K-1} \sum_{c=1}^C F_{m,n,c,d} \cdot x_{i+m,j+n,c}$$

▷ Element-wise multiplication and summation

```
8:     end for
9:   end for
10:  return  $\mathbf{C}$ 
11: end function
```

Algorithm 16 CNN RELU Function for Training and Testing

```
1: function RELU( $\mathbf{C}$ )
2:   Input: Convolved feature map  $\mathbf{C} \in \mathbb{R}^{H' \times W' \times D}$ 
3:   Output: Activated feature map  $\mathbf{A} \in \mathbb{R}^{H' \times W' \times D}$ 
4:   for each element  $(i, j, d)$  in  $\mathbf{C}$  do
5:     Apply ReLU activation function:
```

$$A_{i,j,d} = \max(0, C_{i,j,d})$$

```
6:   end for
7:   return  $\mathbf{A}$ 
8: end function
```

Algorithm 17 CNN MAXPOOLING Function for Training and Testing

```
1: function MAXPOOLING( $\mathbf{A}$ )
2:   Input: Activated feature map  $\mathbf{A} \in \mathbb{R}^{H' \times W' \times D}$ 
3:   Output: Pooled feature map  $\mathbf{P} \in \mathbb{R}^{H'' \times W'' \times D}$ 
4:   for each depth channel  $d \in \{1, \dots, D\}$  do
5:     for each pooling window  $(i, j)$  in  $\mathbf{A}$  do
6:       Apply max pooling operation:
```

$$P_{i,j,d} = \max_{(m,n) \in W} A_{i+m,j+n,d}$$

```
7:     end for
8:   end for
9:   return  $\mathbf{P}$ 
10: end function
```

Algorithm 18 CNN FULLYCONNECTED Function for Training and Testing

- 1: **function** FULLYCONNECTED($\mathbf{P}, \mathbf{W}, \mathbf{b}$)
- 2: **Input:** Pooled feature vector $\mathbf{P} \in \mathbb{R}^d$, weights $\mathbf{W} \in \mathbb{R}^{C \times d}$, biases $\mathbf{b} \in \mathbb{R}^C$
- 3: **Output:** Class probabilities $\hat{\mathbf{y}} \in \mathbb{R}^C$
- 4: **Step 1: Compute class scores (logits)**

$$\mathbf{Z} = \mathbf{WP} + \mathbf{b}$$

where:

- \mathbf{P} is the flattened feature vector from the previous layer.
- \mathbf{W} is the weight matrix connecting features to output neurons.
- \mathbf{b} is the bias vector.
- $\mathbf{Z} \in \mathbb{R}^C$ represents the logits (unnormalized scores for each class).

- 5: **Step 2: Apply softmax activation**

$$\hat{y}_k = \frac{e^{Z_k}}{\sum_{j=1}^C e^{Z_j}}, \quad k = 1, 2, \dots, C$$

where:

- $\hat{\mathbf{y}}$ represents the probability distribution over C classes.
- Each output \hat{y}_k is the normalized score for class k .

- 6: **return** $\hat{\mathbf{y}}$
 - 7: **end function**
-

Algorithm 19 CNN Testing Algorithm for Text Classification

Require: Test sample \mathbf{x} , trained CNN model with filter weights \mathbf{F} , network parameters \mathbf{W}, \mathbf{b}

Ensure: Predicted class label \hat{y}

- 1: **function** PREDICTCNN($\mathbf{x}, \mathbf{F}, \mathbf{W}, \mathbf{b}$)
- 2: **Step 1: Forward Pass**
- 3: $\mathbf{C} \leftarrow \text{CONVOLUTION}(\mathbf{x}, \mathbf{F})$
- 4: $\mathbf{A} \leftarrow \text{ReLU}(\mathbf{C})$
- 5: $\mathbf{P} \leftarrow \text{MAXPOOLING}(\mathbf{A})$
- 6: $\hat{\mathbf{y}} \leftarrow \text{FULLYCONNECTED}(\mathbf{P}, \mathbf{W}, \mathbf{b})$
- 7: **Step 2: Prediction**
- 8: Assign class label:

$$\hat{y} = \arg \max_k \hat{\mathbf{y}}[k]$$

- 9: **return** \hat{y}
 - 10: **end function**
-

Algorithm 11 represents the forward pass through a single convolutional layer in a CNN, involving the following key steps:

- **Convolution:** Each input text representation is processed by convolutional filters to generate an output feature map \mathbf{C} .
- **ReLU Activation:** The Rectified Linear Unit (ReLU) function is applied element-wise to \mathbf{C} to introduce non-linearity.
- **Max Pooling:** Each region in \mathbf{C} is downsampled using max pooling over a fixed-size window to generate the pooled feature map \mathbf{P} .

This sequence of operations is fundamental in CNNs for text classification, as it enables hierarchical feature extraction while reducing computational complexity. The learned features are then passed through fully connected layers for final classification.

6.3 Implementation Details

Data Preprocessing: Proper data preprocessing is essential for enhancing model performance and ensuring stable convergence during training. First, pixel values should be normalized to the range $[0, 1]$ or standardized using zero mean and unit variance to maintain numerical stability. Additionally, data augmentation techniques such as random rotations, flips, and zooms can be applied to improve generalization by artificially increasing the diversity of training samples. For multi-class classification tasks, categorical labels should be encoded as one-hot vectors to allow the model to process class probabilities effectively. Finally, the dataset should be split into training, validation, and testing subsets to enable model training, hyperparameter tuning, and unbiased performance evaluation.

Model Training: Training a Convolutional Neural Network (CNN) involves a structured approach to feature extraction and classification. Convolutional layers with small filters (e.g., 3×3) are employed, followed by max pooling layers to reduce spatial dimensions while preserving important features. The ReLU activation function is applied in convolutional layers to introduce non-linearity, whereas the softmax activation function is used in the output layer to compute class probabilities for classification. Model parameters are optimized using gradient-based optimizers such as Adam or Stochastic Gradient Descent (SGD) with momentum to accelerate convergence. Additionally, hyperparameters including learning rate η , number of epochs T , batch size B , and network depth should be fine-tuned using cross-validation to achieve optimal performance. To mitigate overfitting and improve generalization, regularization techniques such as dropout and batch normalization are incorporated, ensuring stability during training.

Model Evaluation: Evaluating the performance of a Convolutional Neural Network (CNN) involves assessing various metrics to ensure its effectiveness in classification tasks. Standard evaluation metrics such as accuracy, precision, recall, F1-score, and confusion matrices provide insights into the model's predictive capabilities and class-wise performance. Additionally, training and validation loss curves should be monitored to detect signs of overfitting or underfitting, allowing for adjustments in hyperparameters or regularization techniques. To gain a deeper understanding of model predictions, interpretability techniques such as Grad-CAM can be employed to visualize feature maps and highlight the most influential regions in the input data. These evaluation methods collectively help in refining the model and improving its generalization ability.

6.4 Advantages and Limitations

Advantages: Convolutional Neural Networks (CNNs) offer several advantages, making them highly effective for image and spatial data processing. They efficiently capture spatial hierarchies in data by leveraging local receptive fields, enabling hierarchical feature extraction. Compared to fully connected networks, CNNs significantly reduce the number of parameters through weight sharing and local connectivity, making deep architectures computationally feasible. The shared weights in convolutional layers provide translation-invariant feature learning, allowing the model to recognize patterns regardless of their spatial position. CNNs achieve state-of-the-art performance in tasks such as image classification and object detection by effectively extracting meaningful features from raw data. Additionally, pooling operations enhance robustness to small shifts, distortions, and variations in input data, improving the generalization capability of the model.

Limitations: Despite their powerful feature extraction capabilities, Convolutional Neural Networks (CNNs) come with certain limitations. They require large amounts of labeled training data to generalize effectively, making them data-hungry compared to traditional machine learning models. Additionally, CNNs are computationally expensive, especially for deep architectures, often necessitating the use of GPUs or specialized hardware for efficient training. Hyperparameter tuning, including selecting appropriate filter sizes, the number of layers, stride values, and pooling sizes, is complex and highly problem-specific, requiring extensive experimentation. CNNs are also vulnerable to adversarial attacks, where small perturbations in input images can mislead the network into making incorrect predictions. Furthermore, compared to simpler models, CNNs have lower interpretability, making model explainability challenging and limiting their use in applications requiring high transparency.

6.5 Analysis of the CNN Algorithms for Training and Testing

Convolutional Neural Networks (CNNs) are widely used for text classification due to their hierarchical feature extraction capabilities. The training and testing algorithms involve multiple stages, including convolutional operations, activation functions, pooling layers, fully connected layers, and optimization processes. The CNN model captures local spatial dependencies using convolutional filters and learns high-level representations for classification tasks.

Time Complexity Analysis

The time complexity of CNNs is determined by operations in the convolutional, activation, pooling, and fully connected layers. The complexity varies between the training and testing phases.

Training Phase Complexity

The key computational steps in the training phase are as follows:

- **Convolution Operation:** For an input of size $n \times n$ with c channels and f filters of size $k \times k$, the complexity is:

$$O(f \cdot c \cdot k^2 \cdot n^2)$$

- **Activation and Pooling:** The ReLU activation function is applied element-wise with a complexity of $O(n^2)$. Pooling operations (e.g., max pooling) reduce dimensions by a factor of

p , contributing:

$$O(n^2/p^2)$$

- **Fully Connected Layer:** Assuming a hidden layer with h neurons, the complexity is:

$$O(h \cdot n^2)$$

- **Gradient Computation and Weight Updates:** The backpropagation process includes computing gradients for convolutional layers and updating weights, adding an additional:

$$O(f \cdot k^2 \cdot n^2)$$

Thus, the total training complexity is:

$$O(T \cdot N \cdot (f \cdot k^2 \cdot n^2 + h \cdot n^2))$$

where T is the number of epochs and N is the number of training samples.

Testing Phase Complexity

The testing phase only involves a forward pass, skipping weight updates. The complexity remains:

$$O(f \cdot k^2 \cdot n^2 + h \cdot n^2)$$

This is significantly faster than training since gradient computations are omitted.

6.6 Correctness Proof

The correctness of the CNN algorithm is established through its ability to approximate complex functions, extract meaningful features, and minimize classification error using gradient-based optimization. CNNs leverage convolutional layers to identify spatial dependencies, ReLU activation functions to introduce non-linearity, and fully connected layers to map extracted features to class probabilities. Below, we provide a detailed correctness proof.

Proof Outline

1. **Feature Extraction:** Convolutional layers extract hierarchical features using learnable filters that capture spatial structures in the input data. Mathematically, each convolution operation can be expressed as:

$$C_{i,j}^{(l)} = \sum_{m=0}^{k-1} \sum_{n=0}^{k-1} W_{m,n}^{(l)} X_{i+m,j+n} + b^{(l)}$$

where $C_{i,j}^{(l)}$ is the output feature map at position (i, j) for layer l , $W_{m,n}^{(l)}$ represents the convolutional filter weights, $X_{i+m,j+n}$ is the corresponding region in the input, and $b^{(l)}$ is the bias term. This operation enables CNNs to detect patterns such as edges, shapes, and textures, which are crucial for classification.

2. **Non-linearity:** The application of the ReLU (Rectified Linear Unit) activation function introduces non-linearity, allowing CNNs to learn complex representations. The ReLU function is defined as:

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

which ensures that only positive values are propagated forward while setting negative values to zero. This non-linearity enables CNNs to approximate complex decision boundaries and prevents vanishing gradients during training.

3. **Spatial Invariance via Pooling:** Pooling layers, such as max pooling, reduce the dimensionality of feature maps while preserving the most important features. Given a feature map $C^{(l)}$ with window size $p \times p$, max pooling is defined as:

$$P_{i,j}^{(l)} = \max_{(m,n) \in W_{p,p}} C_{i+m,j+n}^{(l)}$$

where $P_{i,j}^{(l)}$ represents the pooled feature map, and $W_{p,p}$ defines the pooling window. Pooling enhances the network's robustness to minor spatial variations in input data.

4. **Optimization Convergence:** The CNN is trained using backpropagation, which adjusts filter weights and biases to minimize a loss function. The gradient descent update rule for a weight parameter W is:

$$W^{(t+1)} = W^{(t)} - \eta \frac{\partial L}{\partial W}$$

where η is the learning rate and $\frac{\partial L}{\partial W}$ is the gradient of the loss function L with respect to W . The iterative update process ensures convergence toward an optimal solution that minimizes classification error.

5. **Classification Consistency:** The fully connected layers aggregate extracted features and map them to class probabilities using the softmax function:

$$P(y = k|x) = \frac{e^{z_k}}{\sum_{j=1}^C e^{z_j}}$$

where z_k is the activation for class k , and C is the total number of classes. The softmax function ensures that output values sum to one, forming a valid probability distribution over the classes. This guarantees that the model makes consistent and interpretable classification predictions.

Mathematical Justification of Convergence

A well-trained CNN minimizes the cross-entropy loss function:

$$L = - \sum_{i=1}^N \sum_{k=1}^C y_{i,k} \log P(y = k|x_i)$$

where $y_{i,k}$ represents the true label of sample i for class k . The gradient-based optimization algorithm ensures that updates decrease this loss iteratively, eventually converging to a local minimum. Under standard assumptions (e.g., convexity in small regions), the Stochastic Gradient Descent (SGD) optimization with momentum leads to a stable and optimal solution [10].

Summary - Correctness Proof

The CNN algorithm is provably correct under the following assumptions:

- **Feature extraction:** The convolutional layers effectively learn hierarchical representations of input data.
- **Non-linearity and invariance:** The combination of ReLU activation and pooling layers enables the model to capture complex patterns while maintaining spatial invariance.
- **Optimization convergence:** The backpropagation algorithm updates network parameters iteratively, leading to convergence to an optimal set of weights that minimizes classification loss.
- **Valid class predictions:** The softmax function ensures that the CNN outputs a valid probability distribution over possible classes, leading to consistent classification results.

Thus, CNNs provide a mathematically sound and empirically validated framework for deep learning-based classification tasks.

6.7 Summary for Convolutional Neural Networks

Convolutional Neural Networks (CNNs) are a fundamental deep learning architecture for processing structured data such as images and text. By leveraging convolutional layers, CNNs automatically extract hierarchical features, enabling them to recognize spatial patterns with minimal preprocessing [18]. The inclusion of ReLU activation functions introduces non-linearity, allowing CNNs to capture complex data relationships [9]. Pooling layers help reduce computational complexity and enhance model generalization by downsampling feature maps, while fully connected layers aggregate extracted features for final classification [17].

CNNs achieve state-of-the-art performance in image classification, object detection, and NLP tasks by efficiently capturing spatial hierarchies [10, 15]. Their training complexity is approximately $O(T \cdot N \cdot d \cdot k^2)$, where T is the number of epochs, N is the number of training samples, d is the number of filters per layer, and $k \times k$ is the kernel size. The testing complexity is significantly lower at $O(M \cdot d \cdot k^2)$, as only forward propagation is performed.

The correctness of CNNs is justified by their ability to approximate complex functions, optimize feature extraction through backpropagation, and maintain classification consistency using softmax-based output layers [10]. Despite their effectiveness, CNNs have limitations, such as their reliance on large labeled datasets, high computational costs, and susceptibility to adversarial attacks [25]. However, advancements in transfer learning, data augmentation, and efficient architectures such as ResNets and EfficientNets continue to enhance CNN performance while mitigating these challenges [13, 26].

Overall, CNNs have established themselves as a cornerstone of deep learning applications, delivering high accuracy, scalability, and adaptability across various domains. When properly tuned and regularized, they achieve superior generalization, making them indispensable for modern artificial intelligence systems.

7 Recurrent Neural Networks (RNNs)

Recurrent Neural Networks (RNNs) are a class of neural networks specifically designed for modeling sequential data, such as text, speech, and time-series signals [14, 1]. Unlike traditional feedforward neural networks, RNNs maintain an internal hidden state that allows them to capture dependencies across different time steps. This recurrent architecture makes them well-suited for various natural language processing (NLP) tasks, including text classification, sentiment analysis, and language modeling [3, 20]. The primary advantage of RNNs lies in their ability to process variable-length input sequences by sequentially updating hidden states, which enables them to learn context-aware representations of sequential data.

Mathematically, at each time step t , an RNN updates its hidden state using the recurrence relation:

$$\mathbf{h}_t = \tanh(\mathbf{W}_x \mathbf{x}_t + \mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{b}_h), \quad (27)$$

where \mathbf{h}_t is the hidden state at time t , \mathbf{x}_t is the input at t , and $\mathbf{W}_x, \mathbf{W}_h$ are the learnable weight matrices with bias \mathbf{b}_h . After processing the input sequence, the final hidden state is passed through a fully connected layer to produce an output probability distribution over class labels:

$$\hat{\mathbf{y}} = \text{Softmax}(\mathbf{W}_y \mathbf{h}_L + \mathbf{b}_y). \quad (28)$$

RNNs are trained using Backpropagation Through Time (BPTT), an extension of standard backpropagation that accounts for dependencies across time steps. Despite their success in sequence modeling, RNNs suffer from vanishing and exploding gradient problems, which limit their ability to capture long-range dependencies [1]. To address this, advanced variants such as Long Short-Term Memory (LSTM) [14] and Gated Recurrent Units (GRUs) [3] have been introduced to improve long-term memory retention and training stability.

7.1 Mathematical Formulations

Given an input sequence $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ of length N , an RNN processes each time step recursively using a hidden state update function:

$$\mathbf{h}_t = \tanh(\mathbf{W}_x \mathbf{x}_t + \mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{b}_h), \quad (29)$$

where:

- \mathbf{h}_t is the hidden state at time step t .
- \mathbf{x}_t is the input vector at time step t .
- \mathbf{W}_x and \mathbf{W}_h are learnable weight matrices.
- \mathbf{b}_h is the bias vector.
- $\tanh(\cdot)$ is the hyperbolic tangent activation function.

The final hidden state \mathbf{h}_L (after processing the entire sequence) is passed through a fully connected layer to obtain the output logits:

$$\mathbf{Z} = \mathbf{W}_y \mathbf{h}_L + \mathbf{b}_y, \quad (30)$$

where \mathbf{W}_y and \mathbf{b}_y are weight and bias parameters of the output layer.

To obtain class probabilities for classification, the softmax function is applied:

$$\hat{y}_k = \frac{e^{Z_k}}{\sum_{j=1}^C e^{Z_j}}, \quad (31)$$

where C is the number of classes.

The loss function used for classification is typically cross-entropy loss:

$$L = - \sum_{i=1}^N \sum_{k=1}^C y_{i,k} \log \hat{y}_{i,k}, \quad (32)$$

where $y_{i,k}$ represents the true label of sample i .

7.2 Description of the Algorithm

The Recurrent Neural Network (RNN) training algorithm [20](#) follows four key steps:

1. Forward Pass:

- The input sequence is processed sequentially, updating the hidden state at each time step.
- The hidden state at time step t is updated using:

$$\mathbf{h}_t = \tanh(\mathbf{W}_x \mathbf{x}_t + \mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{b}_h). \quad (33)$$

- The final hidden state \mathbf{h}_L is passed through a fully connected layer to compute class scores:

$$\mathbf{Z} = \mathbf{W}_y \mathbf{h}_L + \mathbf{b}_y. \quad (34)$$

- The softmax function is applied to obtain class probabilities:

$$\hat{y}_k = \frac{e^{Z_k}}{\sum_{j=1}^C e^{Z_j}}, \quad k \in \{1, 2, \dots, C\}, \quad (35)$$

where C is the number of output classes.

2. Loss Computation:

- The predicted class probabilities \hat{y} are compared with the true labels using the cross-entropy loss function:

$$L = - \sum_{i=1}^N \sum_{k=1}^C y_{i,k} \log \hat{y}_{i,k}. \quad (36)$$

3. Backward Pass (Backpropagation Through Time - BPTT):

- The loss gradients are computed with respect to the model parameters using the chain rule.
- Gradients for the output layer weights:

$$\nabla \mathbf{W}_y, \nabla \mathbf{b}_y = \frac{\partial L}{\partial \mathbf{W}_y}, \frac{\partial L}{\partial \mathbf{b}_y}. \quad (37)$$

- Gradients for the recurrent layer weights:

$$\nabla \mathbf{W}_x, \nabla \mathbf{W}_h, \nabla \mathbf{b}_h = \frac{\partial L}{\partial \mathbf{W}_x}, \frac{\partial L}{\partial \mathbf{W}_h}, \frac{\partial L}{\partial \mathbf{b}_h}. \quad (38)$$

4. Parameter Updates:

- The model parameters are updated using a gradient-based optimizer such as Stochastic Gradient Descent (SGD) or Adam:

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \cdot \nabla \mathbf{W}, \quad (39)$$

where η is the learning rate.

This process is repeated over multiple epochs until the model converges to an optimal solution.

Algorithm 20 RNN Training Algorithm for Text Classification

Require: Training dataset $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$, labels $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$, number of epochs T , learning rate η , network parameters $\mathbf{W}_h, \mathbf{W}_x, \mathbf{W}_y, \mathbf{b}_h, \mathbf{b}_y$

Ensure: Trained RNN model with optimized parameters

```

1: function TRAINRNN( $\mathbf{X}, \mathbf{y}, T, \eta, \mathbf{W}_h, \mathbf{W}_x, \mathbf{W}_y, \mathbf{b}_h, \mathbf{b}_y$ )
2:   for  $t = 1$  to  $T$  do ▷ Epoch loop
3:     for each training sample  $(\mathbf{x}_i, y_i)$  in  $(\mathbf{X}, \mathbf{y})$  do
4:       Initialize hidden state  $\mathbf{h}_0 = \mathbf{0}$ 
5:       Step 1: Forward Pass
6:       for  $t = 1$  to sequence length  $L$  do
7:          $\mathbf{h}_t \leftarrow \text{RECURRENTSTEP}(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{W}_x, \mathbf{W}_h, \mathbf{b}_h)$ 
8:       end for
9:        $\hat{\mathbf{y}} \leftarrow \text{FULLYCONNECTED}(\mathbf{h}_L, \mathbf{W}_y, \mathbf{b}_y)$ 
10:      Step 2: Compute Loss
11:       $L \leftarrow \text{CROSSENTROPYLOSS}(\hat{\mathbf{y}}, y_i)$ 
12:      Step 3: Backward Pass (Gradient Calculation)
13:       $\nabla \mathbf{W}_y, \nabla \mathbf{b}_y \leftarrow \text{COMPUTEGRAIENTS}(L, \mathbf{W}_y, \mathbf{b}_y)$ 
14:       $\nabla \mathbf{W}_x, \nabla \mathbf{W}_h, \nabla \mathbf{b}_h \leftarrow \text{COMPUTERECURRENTGRADIENTS}(L, \mathbf{W}_x, \mathbf{W}_h, \mathbf{b}_h)$ 
15:      Step 4: Parameter Update
16:       $\mathbf{W}_y \leftarrow \mathbf{W}_y - \eta \cdot \nabla \mathbf{W}_y$ 
17:       $\mathbf{b}_y \leftarrow \mathbf{b}_y - \eta \cdot \nabla \mathbf{b}_y$ 
18:       $\mathbf{W}_x \leftarrow \mathbf{W}_x - \eta \cdot \nabla \mathbf{W}_x$ 
19:       $\mathbf{W}_h \leftarrow \mathbf{W}_h - \eta \cdot \nabla \mathbf{W}_h$ 
20:       $\mathbf{b}_h \leftarrow \mathbf{b}_h - \eta \cdot \nabla \mathbf{b}_h$ 
21:    end for
22:  end for
23:  return Trained RNN model with updated parameters
24: end function

```

The RNN testing algorithm 21 evaluates a trained model on unseen data and consists of the following steps:

1. Forward Pass:

- The input sequence is processed sequentially, updating the hidden state at each time step:

$$\mathbf{h}_t = \tanh(\mathbf{W}_x \mathbf{x}_t + \mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{b}_h). \quad (40)$$

- The final hidden state at time step L is passed through a fully connected layer:

$$\mathbf{Z} = \mathbf{W}_y \mathbf{h}_L + \mathbf{b}_y. \quad (41)$$

- The softmax function is applied to compute class probabilities.

2. Prediction:

- The predicted class label is assigned based on the highest probability score:

$$\hat{y} = \arg \max_k \hat{\mathbf{y}}[k]. \quad (42)$$

Unlike training, testing only involves the forward pass and does not require backpropagation or weight updates, making it computationally more efficient.

Algorithm 21 RNN Testing Algorithm for Text Classification

Require: Test sample \mathbf{x} , trained RNN model with parameters $\mathbf{W}_h, \mathbf{W}_x, \mathbf{W}_y, \mathbf{b}_h, \mathbf{b}_y$

Ensure: Predicted class label \hat{y}

```
1: function PREDICTRNN( $\mathbf{x}, \mathbf{W}_h, \mathbf{W}_x, \mathbf{W}_y, \mathbf{b}_h, \mathbf{b}_y$ )
2:   Initialize hidden state  $\mathbf{h}_0 = \mathbf{0}$ 
3:   Step 1: Forward Pass
4:   for  $t = 1$  to sequence length  $L$  do
5:      $\mathbf{h}_t \leftarrow \text{RECURRENTSTEP}(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{W}_x, \mathbf{W}_h, \mathbf{b}_h)$ 
6:   end for
7:    $\hat{\mathbf{y}} \leftarrow \text{FULLYCONNECTED}(\mathbf{h}_L, \mathbf{W}_y, \mathbf{b}_y)$ 
8:   Step 2: Prediction
9:   Assign class label:
```

$$\hat{y} = \arg \max_k \hat{\mathbf{y}}[k]$$

```
10:  return  $\hat{y}$ 
11: end function
```

7.3 Implementation Details

Data Preprocessing: Proper data preprocessing is essential for ensuring the effective training of a Recurrent Neural Network (RNN) for text classification. First, text data is tokenized and converted into numerical representations using word embeddings such as Word2Vec, GloVe, or pre-trained embeddings from transformer models. This step captures semantic relationships between words and enhances the network’s ability to understand language structures. Next, input sequences are padded or truncated to a fixed length to enable efficient batch processing and maintain uniform input dimensions. Categorical labels are encoded as one-hot vectors, allowing the model to handle multi-class classification tasks. Additionally, input data is normalized to improve training stability by ensuring consistent numerical scales across features. Finally, the dataset is split into training, validation, and testing subsets to facilitate model training, hyperparameter tuning, and unbiased performance evaluation.

Model Processing: The processing of the Recurrent Neural Network (RNN) involves initializing and updating model parameters during training. First, the model parameters, including weight matrices \mathbf{W}_x , \mathbf{W}_h , \mathbf{W}_y and bias vectors \mathbf{b}_h , \mathbf{b}_y , are initialized. During training, the input sequences propagate through the recurrent layers, where each time step updates the hidden state based on the current input and the previous hidden state. The model then computes gradients using Backpropagation Through Time (BPTT), an extension of backpropagation that accounts for dependencies across time steps in sequential data. Finally, the model updates parameters using an optimization algorithm such as Stochastic Gradient Descent (SGD) or Adam, ensuring the network learns meaningful representations from the input data.

Model Evaluation: Evaluating the performance of a Recurrent Neural Network (RNN) involves assessing various metrics to ensure its effectiveness in classification tasks. The primary evaluation metrics include accuracy, precision, recall, and F1-score, which provide insights into the model's predictive capabilities and class-wise performance. Additionally, training and validation loss curves should be monitored to detect potential overfitting or underfitting, allowing for adjustments in hyperparameters or regularization techniques. To enhance model interpretability, attention-based methods can be utilized to visualize word importance, offering insights into how the model makes predictions and which parts of the input sequences contribute most significantly to the classification decision.

7.4 Advantages and Limitations

Advantages: Recurrent Neural Networks (RNNs) offer several advantages, making them highly effective for sequential data processing tasks such as text classification and natural language processing (NLP). One of the key strengths of RNNs is their ability to capture sequential dependencies, enabling the model to learn context across time steps. Additionally, RNNs can process variable-length input sequences, making them well-suited for text and speech applications where input sizes may vary. Parameter sharing across time steps significantly reduces memory complexity compared to fully connected networks, improving computational efficiency. Furthermore, RNNs can effectively handle both short and long sequences when enhanced with techniques such as bidirectional RNNs or attention mechanisms, which allow the model to learn dependencies from both past and future context, leading to improved performance in complex NLP tasks.

Limitations: Despite their strengths, Recurrent Neural Networks (RNNs) exhibit several limitations that affect their performance and scalability. One major drawback is the vanishing and exploding gradient problem, which makes training deep RNNs challenging as gradients tend to decay or grow exponentially over long sequences [1]. As a result, RNNs struggle with long-range dependencies, where information from earlier time steps is lost due to the exponential decay of gradients. Additionally, RNNs are computationally expensive because their sequential nature prevents efficient parallelization, unlike Convolutional Neural Networks (CNNs), which can process inputs simultaneously. Another challenge is interpretability, as RNNs are less transparent compared to simpler models, making it difficult to explain their decision-making process. Lastly, training requires careful tuning of hyperparameters such as learning rate, sequence length, and hidden state size, which can significantly impact model performance and generalization.

7.5 Analysis of the RNN Algorithms for Training and Testing

The training and testing of an RNN involve sequential data processing, hidden state updates, and weight optimization. The algorithm operates through the following key stages:

Training Phase

The training process of an RNN consists of:

- Forward Pass: The input sequence is processed sequentially, with each time step updating the hidden state using:

$$\mathbf{h}_t = \tanh(\mathbf{W}_x \mathbf{x}_t + \mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{b}_h) \quad (43)$$

- Loss Computation: The predicted output at the final time step is compared to the true label using the cross-entropy loss function:

$$L = - \sum_{k=1}^C y_k \log(\hat{y}_k), \quad (44)$$

where \hat{y}_k represents the softmax probability of class k .

- Backward Pass (Backpropagation Through Time - BPTT): The gradients of the loss function are computed with respect to the parameters:

$$\frac{\partial L}{\partial \mathbf{W}_x}, \quad \frac{\partial L}{\partial \mathbf{W}_h}, \quad \frac{\partial L}{\partial \mathbf{W}_y} \quad (45)$$

These gradients are accumulated through time steps, making the process computationally expensive.

- Parameter Updates: The parameters are updated using an optimization algorithm such as SGD or Adam:

$$\theta \leftarrow \theta - \eta \frac{\partial L}{\partial \theta} \quad (46)$$

where θ represents the learnable parameters.

Testing Phase

During testing, the trained RNN processes an input sequence without updating the parameters:

- The input sequence propagates through the network, updating the hidden states.
- The final hidden state is passed through a fully connected layer to generate class scores.
- The predicted class label is determined using:

$$\hat{y} = \arg \max_k \hat{\mathbf{y}}_k. \quad (47)$$

7.6 Time Complexity Analysis

The time complexity of RNNs is determined by the number of time steps L , the input size n , and the hidden state size h . The major computations are:

- Forward Pass Complexity: Each time step involves matrix multiplications:

$$O(L \cdot (nh + h^2)) \quad (48)$$

where nh accounts for input-to-hidden computations and h^2 for hidden-to-hidden transitions.

- Backward Pass Complexity: The backpropagation through time (BPTT) requires computing gradients for each time step:

$$O(L \cdot (nh + h^2)) \quad (49)$$

making it as expensive as the forward pass.

- Total Training Complexity: For N training samples over T epochs:

$$O(T \cdot N \cdot L \cdot (nh + h^2)). \quad (50)$$

- Testing Complexity: Since testing only involves a forward pass:

$$O(L \cdot (nh + h^2)). \quad (51)$$

7.7 Correctness Proof

The correctness of the RNN algorithm is based on its ability to capture sequential dependencies and optimize classification performance.

Proof Outline

1. Sequential Processing: Each input \mathbf{x}_t is mapped to a hidden state \mathbf{h}_t that retains historical context. The recurrence relation:

$$\mathbf{h}_t = \tanh(\mathbf{W}_x \mathbf{x}_t + \mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{b}_h) \quad (52)$$

ensures that information from previous time steps influences future predictions.

2. Gradient-Based Optimization: The backpropagation through time (BPTT) algorithm ensures that the gradient $\frac{\partial L}{\partial \mathbf{W}}$ is propagated backward through the network. The loss function:

$$L = - \sum_{k=1}^C y_k \log(\hat{y}_k) \quad (53)$$

is minimized through iterative gradient descent updates.

3. Convergence: Assuming the loss function is differentiable and the learning rate η is properly tuned, gradient descent ensures convergence to a local minimum.
4. Classification Consistency: The softmax output:

$$P(y = k|x) = \frac{e^{z_k}}{\sum_{j=1}^C e^{z_j}} \quad (54)$$

guarantees that the highest probability class is selected as the final prediction.

Summary - Correctness Proof

The RNN algorithm is theoretically sound under the following assumptions:

- Sequential dependency is preserved through recurrent connections.
- Gradient-based optimization leads to convergence given a well-chosen learning rate.
- Classification outputs are valid probability distributions due to the softmax function.
- Empirical validation confirms that RNNs can effectively model sequential dependencies in text classification tasks [20].

7.8 Summary of Recurrent Neural Networks

Recurrent Neural Networks (RNNs) are powerful models for sequential data processing, offering the ability to capture temporal dependencies through recurrent connections. Their architecture enables them to effectively model structured sequences, making them particularly useful for natural language processing and time-series forecasting. The training complexity of an RNN is approximately:

$$O(T \cdot N \cdot L \cdot d^2),$$

where:

- T is the number of training epochs,
- N is the number of training samples,
- L is the sequence length,
- d is the hidden state size.

The testing complexity is:

$$O(N_{\text{test}} \cdot L \cdot d^2),$$

where N_{test} is the number of test samples.

While RNNs are highly effective for capturing short-term dependencies, they struggle with long-range dependencies due to the vanishing gradient problem. This issue can be mitigated by using LSTM and GRU architectures, which introduce gating mechanisms to selectively retain important information over long sequences. Additionally, RNNs are computationally expensive due to their sequential nature, limiting parallelization compared to feedforward networks and convolutional neural networks (CNNs) [12].

Despite these challenges, RNNs remain a fundamental building block in deep learning, particularly for NLP applications such as text classification, speech recognition, and machine translation. With advancements in attention mechanisms and transformer-based architectures, modern deep learning models now integrate RNNs with self-attention to achieve state-of-the-art performance across various sequence-based tasks.

8 Recent Advances in Neural Networks and Deep Learning

Recent developments in neural networks and deep learning have transformed the landscape of machine learning, enabling more sophisticated and scalable models for diverse applications. Frameworks such as PyTorch have lowered the barrier to entry for both researchers and practitioners, providing tools for rapid prototyping, efficient computation, and advanced optimization on modern hardware. Meanwhile, Hugging Face has revolutionized Natural Language Processing (NLP) by offering easy-to-use, pre-trained transformer models that can be quickly fine-tuned for various tasks. These advancements, combining the power of flexible deep learning libraries with large-scale pre-trained models, have catalyzed breakthroughs in domains such as text classification, computer vision, and beyond.

8.1 PyTorch: Revolutionizing Deep Learning Frameworks

PyTorch has emerged as one of the most popular deep learning frameworks, particularly within the research community. Its flexibility, dynamic computation graph, and seamless integration with Python make it highly accessible for developing and experimenting with deep neural networks. Recent advancements in PyTorch have centered on the following areas:

- **Autograd and Tensor Operations:** PyTorch's automatic differentiation (**Autograd**) and tensor operations are continuously optimized for performance on both CPUs and GPUs. In recent versions, improved support for mixed precision training (float16 operations) has significantly accelerated model training without compromising accuracy [4].
- **TorchScript and JIT Compilation:** The introduction of **TorchScript** has enabled exporting and deploying models to production with enhanced speed and efficiency. The Just-In-Time (JIT) compiler optimizes PyTorch code, leading to better inference performance in real-world applications.
- **Distributed Training:** PyTorch provides robust support for multi-GPU and distributed training through `torch.distributed`, facilitating both data parallelism and model parallelism. This has allowed scalable training across multiple nodes, reducing total training time on large datasets.
- **Interoperability with TensorFlow:** PyTorch/TensorFlow interoperability has made it simpler to convert models between these frameworks, encouraging collaboration and unifying workflows across research and production environments.

Thanks to its dynamic graph execution and user-friendly APIs, PyTorch has become a standard for deep learning research and deployment in various industries.

8.2 PyTorch Example: Simple Neural Network for Classification

Below is a Python implementation of a simple neural network using PyTorch for binary classification. The code leverages the `torch.nn.Module` class to define a network with one hidden layer:

```
import torch
import torch.nn as nn
import torch.optim as optim
```

```

from torch.utils.data import DataLoader, TensorDataset

# Sample dataset (features and labels)
X = torch.randn(100, 10) # 100 samples, 10 features
y = torch.randint(0, 2, (100, 1), dtype=torch.float32) # Binary labels

# Create a DataLoader for batch processing
dataset = TensorDataset(X, y)
train_loader = DataLoader(dataset, batch_size=16, shuffle=True)

# Define a simple neural network model
class SimpleNN(nn.Module):
    def __init__(self):
        super(SimpleNN, self).__init__()
        self.fc1 = nn.Linear(10, 32) # Input layer (10 features) -> Hidden layer (32 units)
        self.fc2 = nn.Linear(32, 1) # Hidden layer (32 units) -> Output layer (1 unit for binary)
        self.sigmoid = nn.Sigmoid() # Sigmoid activation for binary output

    def forward(self, x):
        x = torch.relu(self.fc1(x)) # Apply ReLU activation to the hidden layer
        x = self.fc2(x) # Apply output layer
        return self.sigmoid(x) # Sigmoid activation for final output

# Initialize the model, loss function, and optimizer
model = SimpleNN()
criterion = nn.BCELoss() # Binary Cross-Entropy Loss for binary classification
optimizer = optim.Adam(model.parameters(), lr=0.001)

# Training loop
num_epochs = 20
for epoch in range(num_epochs):
    for batch_idx, (data, target) in enumerate(train_loader):
        # Zero the gradients
        optimizer.zero_grad()

        # Forward pass
        output = model(data)

        # Compute the loss
        loss = criterion(output, target)

        # Backward pass and optimization
        loss.backward()
        optimizer.step()

    # Print progress every epoch

```

```
print(f'Epoch [{epoch+1}/{num_epochs}], Loss: {loss.item():.4f}')
```

This example illustrates several important aspects of PyTorch:

- **Dataset Handling:** Synthetic feature and label data are created using `torch.randn` and `torch.randint`, respectively. A `DataLoader` is used for batch-based training.
- **Model Definition:** The `SimpleNN` class represents a basic neural network with a single hidden layer, employing ReLU and Sigmoid activations.
- **Loss and Optimization:** Binary Cross-Entropy loss (`BCELoss`) and the Adam optimizer are used to update the model's parameters.
- **Training Process:** Over the course of multiple epochs, the model undergoes forward passes, loss computations, and backward passes, adjusting parameters to minimize the loss function.

8.3 Hugging Face: Democratizing NLP with Pre-Trained Models

Hugging Face has gained widespread recognition in the field of Natural Language Processing (NLP) by offering open-source, pre-trained models for a variety of tasks. Their transformer-based architectures, such as BERT, GPT-2, and T5, have reshaped state-of-the-art NLP research. Recent efforts in Hugging Face and the `transformers` library focus on user-friendly interfaces and extensive model repositories [7].

- **Model Hub:** The Hugging Face Model Hub hosts a vast range of pre-trained models for tasks including text classification, summarization, question answering, translation, and more [6].
- **Adapters:** Innovative methods, such as adapters, enable efficient fine-tuning of large transformers using fewer parameters than standard approaches. This improves model deployment in environments with limited computational resources [22].
- **Pipeline Abstractions:** Hugging Face provides high-level pipeline abstractions that simplify workflows for common NLP tasks. Users can quickly leverage text classification or summarization with minimal code [27].
- **Datasets Repository:** Hugging Face also offers a repository of publicly available datasets, facilitating data sharing and streamlined model training [19].
- **Inference API and AutoNLP:** The Inference API allows large-scale model deployment, while AutoNLP assists users in fine-tuning models on custom datasets without extensive architecture knowledge [5].

Through Hugging Face's integrated ecosystem, advanced NLP models are more readily accessible, benefiting both research and industry [8].

8.4 The Intersection of PyTorch and Hugging Face

PyTorch's flexibility, combined with Hugging Face's pre-trained models, has accelerated progress in natural language processing, computer vision, and emerging multimodal applications:

- **Seamless Integration:** Models from Hugging Face can be easily imported into PyTorch, leveraging the framework's dynamic graph for efficient experimentation and model refinement.

- **Faster Prototyping:** Researchers and developers can prototype and deploy cutting-edge NLP solutions quickly by combining pre-trained transformers with PyTorch’s straightforward model-building tools.

Overall, the synergy between PyTorch and Hugging Face fosters an environment where deep learning innovation is accessible to a broader audience, supporting both scientific research and commercial product development.

8.5 Summary - Recent Advances in NN and DL

In this section, we explored two leading contributions to deep learning innovation: PyTorch and Hugging Face. PyTorch’s dynamic computation graph and user-friendly APIs have become a mainstay in research and industry, supporting distributed training, mixed-precision arithmetic, and seamless conversion between frameworks. Hugging Face, with its extensive model hub and pipeline abstractions, has democratized access to advanced NLP capabilities, fostering rapid development and deployment of state-of-the-art models. Together, these frameworks exemplify how open-source collaboration and comprehensive tooling can accelerate progress in deep learning, making powerful models and techniques accessible to a broader audience of developers, researchers, and organizations.

9 Neural Networks and Deep Learning using Python Packages

Deep learning has become the cornerstone of modern machine learning, driving innovations in fields such as computer vision, natural language processing, and time-series forecasting. Python provides an extensive range of packages that facilitate the development, training, and deployment of various neural network architectures, from classical Multi-Layer Perceptrons (MLPs) to more advanced models like Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs), Long Short-Term Memory (LSTM) networks, and Radial Basis Function (RBF) networks.

9.1 Packages for Neural Networks and Deep Learning

The following Python packages are broadly categorized based on their suitability for different neural network architectures and use cases:

Foundational Deep Learning Frameworks

- **PyTorch:** Known for its dynamic computation graph and **Autograd** functionality, PyTorch is widely used for research and rapid prototyping. It excels at implementing diverse models, including RBF networks, MLPs, LSTMs, CNNs, and RNNs.
- **TensorFlow / Keras:** A popular ecosystem providing high-level APIs (**Keras**) and low-level control (core **TensorFlow**). Ideal for building and deploying production-scale models, from simple MLPs to large-scale CNNs and RNNs.

Neural Network APIs and Utilities

- **Scikit-learn:** Offers straightforward interfaces for basic neural networks (e.g., **MLPClassifier**) and integrates well with other traditional machine learning algorithms. Although not specialized for deep networks, it is helpful for quick experimentation and smaller datasets.
- **Ecosystem Extensions (e.g., Lightning, Ignite):** Tools like PyTorch Lightning and Ignite abstract away boilerplate code, making it easier to implement and train advanced architectures such as LSTMs, CNNs, and RBF-based networks with minimal overhead.

Specialized Libraries for Model Variants

- **PyTorch Geometric (PyG):** Extends PyTorch to handle graph-structured data, useful for combining neural networks with graph machine learning—particularly relevant for advanced RNN or CNN adaptations in node classification and link prediction.
- **Hugging Face Transformers:** Focuses on transformer-based architectures for tasks like language modeling and text classification. Can complement RNN or LSTM solutions by providing state-of-the-art pretrained models.
- **AutoML and Hyperparameter Tuning (e.g., AutoKeras, Optuna):** Simplifies experimentation by automating architecture search and hyperparameter tuning, supporting various neural network architectures from MLP to CNN to LSTM.

9.2 Application of Packages to Specific Neural Network Architectures

The packages mentioned above are best suited for implementing specific network types and tasks:

- **RBF Networks:** - PyTorch and TensorFlow both allow flexible custom layers for RBF neurons. These networks can be efficiently prototyped using core APIs or high-level frameworks.
- **Multi-Layer Perceptrons (MLPs):** - Scikit-learn provides a quick entry point for MLPs, while PyTorch and TensorFlow offer more flexibility and performance for deeper, more complex feedforward architectures.
- **Convolutional Neural Networks (CNNs):** - PyTorch and TensorFlow have built-in support for convolutional layers, enabling tasks like image classification, object detection, and segmentation.
- **Recurrent Neural Networks (RNNs) and LSTM:** - PyTorch and TensorFlow natively support RNN and LSTM layers for sequence modeling tasks such as text classification, language modeling, and time-series forecasting.
- **Hybrid Architectures:** - Combining RBF layers with CNNs or LSTMs is possible using flexible frameworks like PyTorch Lightning, which allow custom architectures that incorporate multiple types of neural network layers.

9.3 Advantages and Limitations of Python Deep Learning Libraries

Advantages:

- **Rich Ecosystem:** Python's deep learning packages offer comprehensive support for various architectures (MLP, CNN, RNN, LSTM, RBF), alongside extensive tooling for data preprocessing, visualization, and deployment.
- **Ease of Use:** High-level APIs (e.g., Keras, Lightning) simplify network definition and training loops, lowering the barrier for newcomers.
- **Performance and Scalability:** Libraries such as PyTorch and TensorFlow leverage GPUs and distributed training to handle large-scale problems efficiently.

Limitations:

- **Complexity in Advanced Models:** Implementing hybrid or highly specialized architectures (e.g., RBF-CNN or LSTM-GAN) can be challenging and requires extensive framework knowledge.
- **Hardware Requirements:** Large-scale training often demands considerable computational resources (GPUs, TPUs), which may not be easily accessible to all users.

9.4 Summary - Neural Networks and Deep Learning using Python Packages

Python's flourishing deep learning ecosystem enables researchers and practitioners to design and experiment with a diverse range of neural network architectures—including RBF networks, MLPs, CNNs, LSTMs, and RNNs. By selecting the appropriate library or combination of libraries, users

can readily build, train, and deploy models for tasks in computer vision, natural language processing, and beyond. This versatility and community-driven development ensure continued innovation and support for state-of-the-art deep learning solutions.

10 Summary

Neural networks and deep learning form the backbone of contemporary AI research and applications. This document surveyed five key neural network architectures:

- **RBF Networks**—which employ radial basis functions for localized, robust feature representations.
- **Multi-Layer Perceptrons (MLPs)**—the foundational feedforward networks capable of modeling non-linear mappings across a range of problems.
- **Long Short-Term Memory (LSTM)**—a specialized form of RNN designed to maintain long-range dependencies in sequential data.
- **Convolutional Neural Networks (CNNs)**—capable of hierarchical feature extraction, originally devised for image processing but equally adept at structured or sequential data.
- **Recurrent Neural Networks (RNNs)**—enabling sequential modeling through hidden states shared across time steps, though subject to vanishing/exploding gradients without modifications such as LSTM.

Each architecture presents distinct advantages, from localized modeling (RBF) to translation-invariant feature detection (CNN) or memory-based sequence handling (LSTM/RNN). Although their inner mechanisms differ, all rely on iterative optimization—gradient-based backpropagation—to adjust parameters for minimal loss on training data. In practice, the choice of architecture often depends on data characteristics: CNNs excel at handling grid-like structures (e.g., images), RNNs and LSTMs are suited for sequential or temporal patterns, while MLPs and RBF networks can be effective on moderately sized tabular datasets or when simpler interpretability is desired.

Advances in computational hardware and Python-based deep learning frameworks have greatly accelerated the design and experimentation of these architectures, contributing to rapid progress in fields like computer vision, natural language processing, and robotics. As the complexity and scope of AI challenges expand, developers and researchers will continue to push the frontiers of these fundamental neural network designs, integrating innovations such as attention mechanisms, hybrid ensembles, and large-scale unsupervised pre-training. By understanding the fundamental theories, training processes, and implementation details of RBF networks, MLPs, LSTMs, CNNs, and RNNs, practitioners are well-equipped to build solutions addressing both current and future challenges in deep learning.

References

- [1] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. “Learning long-term dependencies with gradient descent is difficult”. In: *IEEE Transactions on Neural Networks*. Vol. 5. 2. IEEE. 1994, pp. 157–166.
- [2] David S Broomhead and David Lowe. “Multivariate functional interpolation and adaptive networks”. In: *Complex Systems* 2.3 (1988), pp. 321–355.
- [3] Kyunghyun Cho et al. “Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation”. In: *Empirical Methods in Natural Language Processing (EMNLP)*. 2014.
- [4] PyTorch Contributors. *PyTorch 1.10 Release*. <https://pytorch.org/blog/pytorch-1.10-released/>. Accessed 2021-10-30. 2021.
- [5] Hugging Face. *AutoNLP by Hugging Face*. <https://huggingface.co/autonlp>. Accessed 2021-08-01. 2021.
- [6] Hugging Face. *Hugging Face Model Hub*. <https://huggingface.co/models>. Accessed 2022-05-10. 2022.
- [7] Hugging Face. *Hugging Face Transformers*. <https://github.com/huggingface/transformers>. Accessed 2021-10-30. 2021.
- [8] Hugging Face. *Hugging Face: Democratizing NLP, CV, and beyond*. <https://huggingface.co/>. Accessed 2023-01-15. 2023.
- [9] Xavier Glorot and Yoshua Bengio. “Understanding the difficulty of training deep feedforward neural networks”. In: *International Conference on Artificial Intelligence and Statistics*. 2011.
- [10] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. URL: <https://doi.org/10.1117/12.2664346>.
- [11] Alex Graves. *Supervised Sequence Labelling with Recurrent Neural Networks*. Springer, 2012.
- [12] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. “Speech recognition with deep recurrent neural networks”. In: *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE. 2013, pp. 6645–6649.
- [13] Kaiming He et al. “Deep Residual Learning for Image Recognition”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [14] Sepp Hochreiter and Jrgen Schmidhuber. “Long Short-Term Memory”. In: *Neural Computing* 9 (8 1997), pp. 1735–1780.
- [15] Yoon Kim. “Convolutional Neural Networks for Sentence Classification”. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2014.
- [16] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *Proceedings of the International Conference on Learning Representations (ICLR)*. 2014.
- [17] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. “ImageNet Classification with Deep Convolutional Neural Networks”. In: (2012), pp. 1097–1105.
- [18] Yann LeCun et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.
- [19] Quentin Lhoest et al. *Datasets: A Community Library for Natural Language Processing*. <https://github.com/huggingface/datasets>. Accessed 2021-07-20. 2021.
- [20] Tomas Mikolov et al. “Recurrent neural network based language model”. In: *Eleventh annual conference of the international speech communication association*. 2010, pp. 1045–1048.
- [21] Jun Park and Irving W Sandberg. “Universal approximation using radial-basis-function networks”. In: *Neural Computation* 3.2 (1991), pp. 246–257.
- [22] Jonas Pfeiffer et al. “AdapterFusion: Non-Destructive Task Composition for Transfer Learning”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics. 2020, pp. 4871–4886.
- [23] D.E. Rumelhart, G.E. Hinton, and R.J. Williams. “Learning Representations by Back-Propagating Errors”. In: *Nature* 323.6088 (1986), pp. 533–536.

- [24] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. “Learning representations by back-propagating errors”. In: *Nature* 323.6088 (1986), pp. 533–536.
- [25] Christian Szegedy et al. “Intriguing properties of neural networks”. In: *Proceedings of the International Conference on Learning Representations (ICLR)*. 2013.
- [26] Mingxing Tan and Quoc V. Le. “EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks”. In: *Proceedings of the International Conference on Machine Learning (ICML)*. 2019.
- [27] Thomas Wolf et al. *Transformers: State-of-the-Art Natural Language Processing*. <https://arxiv.org/abs/1910.03771>. 2020. arXiv: [arXiv:1910.03771](https://arxiv.org/abs/1910.03771) [cs.CL].