| Roll No : A-25 | Name : SHRADDHA SAWANT |
|---|---|
| Class : TE (AI &DS) | Batch : A-2 |
| Date of Performance : | Date of Submission : 16/4/2025 |
| Grade : | |

# EXPERIMENT NO - 08

## AIM : Application of studied algorithms as case study.

**THEORY :**

**Case Study 1: Linear Regression**

Linear regression is a foundational supervised learning algorithm that models the linear relationship between a dependent variable and one or more independent variables. The primary goal of linear regression is to find the best-fitting line (in the case of simple linear regression with one independent variable) or hyperplane (in the case of multiple linear regression with several independent variables) that describes this relationship, allowing for the prediction of the dependent variable based on the values of the independent variables.

Several publicly available datasets are suitable for demonstrating the application of linear regression.

- **Boston Housing Dataset:** This dataset, originally part of the UCI Machine Learning Repository, contains information about housing values in the suburbs of Boston, collected by the U.S. Census Service.[1] It comprises 506 instances, each described by 13 continuous or binary features, such as the per capita crime rate by town, the proportion of residential land zoned for lots over 25,000 sq. ft., the proportion of non-retail business acres per town, and a binary variable indicating if the property is near the Charles River.[1] Linear regression can be applied to this dataset to predict the median value of owner-occupied homes based on these features.[1] Before training the model, it is prudent to check for and handle any missing values in the dataset.[1] While basic linear regression is less sensitive to the magnitude of features, feature scaling can be beneficial in some implementations. Training a linear regression model on this dataset typically involves using libraries like scikit-learn in Python, where the dataset can even be loaded directly.[1] To evaluate the model's performance on unseen data, the dataset is often split into training and testing sets using functions like train_test_split.[1] Common evaluation metrics for regression tasks, including R-squared, Mean Squared Error (MSE), and Mean Absolute Error (MAE), are used to assess the accuracy and goodness of

fit of the model.[7] The strong correlation observed between features like the average number of rooms ('RM') and the lower status of the population ('LSTAT') with the median house value ('MEDV') suggests a relationship where larger homes in more affluent areas tend to command higher prices.[1]

- **Medical Insurance Costs Dataset:** Inspired by the book "Machine Learning with R," this dataset contains 1338 rows of data with medical information and costs billed by health insurance companies.[10] The features include age, gender, BMI, number of children, smoker status, region, and the target variable, insurance charges.[11] Linear regression can be used to predict medical insurance costs based on these attributes.[8] Preprocessing this dataset involves handling categorical features such as gender and region, often through techniques like one-hot encoding to convert them into a numerical format suitable for regression models.[8] Feature scaling can also be important to ensure that features with larger ranges do not disproportionately influence the model. Training a multiple linear regression model on this data can be done using Python libraries such as scikit-learn.[8] The performance of the model is typically evaluated using R-squared, MSE, and MAE.[8] Analysis of this dataset often reveals the significant impact of the 'smoker' feature on insurance charges, reflecting a recognized relationship by insurance providers.[9] Age and BMI also tend to show positive correlations with insurance costs, indicating increased healthcare needs with age and higher body mass index.[9]

- **WHO Statistics on Life Expectancy Dataset:** Compiled by the World Health Organization and the United Nations, this dataset tracks factors that affect life expectancy across 179 countries from 2000 to 2015.[21] It contains 2938 rows and 22 columns, including country, year, developing status, adult mortality, life expectancy, infant deaths, alcohol consumption per capita, country's expenditure on health, immunization coverage, BMI, deaths under 5-years-old, deaths due to HIV/AIDS, GDP, population, body condition, income information, and education.[10] Multiple linear regression can be employed to predict life expectancy based on these diverse health and socio-economic indicators.[23] Preprocessing involves handling missing values, which are present in some columns.[22] Categorical features like 'developing status' need to be encoded numerically. Given the wide range of scales among the features, feature scaling is crucial. Potential multicollinearity among features should also be considered. The model can be trained using Python libraries [25], and its performance can be evaluated using R-squared, MSE, and potentially other metrics relevant to time series analysis if the 'year' feature is incorporated in a time-dependent manner. Higher infant survival rates and lower adult mortality rates likely contribute directly to increased life expectancy.[23] Factors such as access to healthcare (indicated by health expenditure and immunization coverage) and socio-economic status (GDP, income composition, schooling) are also expected to show positive associations with life expectancy, reflecting their broader influence on health outcomes.[23]

The general steps for applying linear regression involve first loading the dataset and conducting exploratory data analysis to understand its characteristics. This may include visualizing

relationships between variables and identifying potential issues like outliers or missing data. Feature engineering might be necessary to create new features or transform existing ones to improve the model's performance. After splitting the data into training and testing sets, a linear regression model is selected and trained on the training data. Finally, the model's performance is evaluated on the testing data using appropriate metrics.

Evaluation metrics for linear regression provide insights into the model's accuracy and how well it fits the data. R-squared, also known as the coefficient of determination, represents the proportion of the variance in the dependent variable that is predictable from the independent variables. It ranges from 0 to 1, with higher values indicating a better fit. Mean Squared Error (MSE) measures the average of the squares of the errors, where an error is the difference between the predicted and actual values. Lower MSE values indicate better performance. Mean Absolute Error (MAE) measures the average of the absolute differences between the predicted and actual values, providing a more robust measure to outliers compared to MSE. The choice of metric depends on the specific goals of the analysis and the characteristics of the data.

## 3. Case Study 2: Logistic Regression

Logistic regression is a statistical model that, despite its name, is primarily used for binary classification problems. It predicts the probability of a binary outcome (e.g., yes/no, true/false) occurring based on one or more predictor variables. Unlike linear regression, which predicts a continuous outcome, logistic regression uses a sigmoid function (also known as the logistic function) to map the linear combination of input features to a probability between 0 and 1.[29]

Several publicly available datasets can be used to demonstrate logistic regression.

- **Iris Dataset (Binary Classification):** The Iris dataset, a classic in machine learning, contains measurements of sepal length, sepal width, petal length, and petal width for 150 iris flowers, belonging to three different species: Iris setosa, Iris versicolor, and Iris virginica.[32] For binary classification using logistic regression, the dataset can be adapted by selecting only two of the three species, for instance, Iris setosa and Iris versicolor.[33] The goal is to classify these two species based on their sepal and petal measurements.[29] Preprocessing may involve feature scaling to ensure that no single feature dominates the model.[33] Training a logistic regression model on this binary subset of the Iris dataset can be done using scikit-learn in Python.[29] The performance of the binary classifier is typically evaluated using metrics such as Accuracy (the proportion of correctly classified instances) [29], Precision (the ratio of correctly predicted positive observations to the total predicted positives), Recall (the ratio of correctly predicted positive observations to the all observations in actual class), F1-score (the harmonic mean of Precision and Recall) [29], and the Area Under the Receiver Operating Characteristic curve (ROC AUC), which measures the ability of the classifier to distinguish between the two classes.[29] A Confusion Matrix, which visualizes the counts of true positives, true negatives, false positives, and false

negatives, is also a valuable evaluation tool.[31] The effectiveness of petal length and petal width in distinguishing between different iris species, even in a binary classification scenario, suggests a fundamental biological difference reflected in these measurements.[33]

- **Titanic Dataset:** This dataset provides information about passengers on the Titanic, including features like passenger class, sex, age, number of siblings/spouses aboard, number of parents/children aboard, fare, and port of embarkation, with the target variable indicating whether a passenger survived the sinking.[40] Logistic regression can be used to predict the survival of passengers based on these features.[40] Preprocessing this dataset is crucial and involves handling missing values, particularly in the 'Age', 'Cabin', and 'Embarked' columns.[40] Categorical features like 'Sex' and 'Embarked' need to be encoded into numerical representations, for example, using one-hot encoding or label encoding.[41] The logistic regression model can then be trained using scikit-learn.[41] Evaluation metrics commonly used include Accuracy, Precision, Recall, F1-score, ROC AUC, and the Confusion Matrix.[40] The higher survival rate observed among female passengers and those in higher passenger classes suggests a relationship influenced by factors such as the "women and children first" protocol and better access to resources for first-class passengers.[40]
- **Predicting Customer Ad Clicks Dataset:** This dataset contains features such as the daily time spent on a website, age, area income, daily internet usage, and gender, with the target variable indicating whether a user clicked on an online advertisement.[30] Logistic regression can be applied to predict the likelihood of a user clicking on an ad.[30] Preprocessing might involve feature scaling depending on the implementation. Training the logistic regression model follows standard procedures. The model's performance is assessed using metrics like Accuracy, Precision, Recall, F1-score, ROC AUC, and the Confusion Matrix. This application demonstrates how online behavior and demographics can be used to predict user engagement with advertisements, potentially reflecting a connection between internet usage, income, and purchasing intent.[30]

The general steps for applying logistic regression include loading and exploring the data, handling categorical features by encoding them numerically, and potentially scaling numerical features. The dataset is then split into training and testing sets. A logistic regression model is instantiated and trained on the training data. After training, the model is used to make predictions on the test data, and its performance is evaluated using appropriate classification metrics.

Evaluation metrics for binary classification provide insights into the model's ability to correctly classify instances. Accuracy measures the overall correctness of the model. Precision quantifies how often the model correctly predicts the positive class. Recall measures the model's ability to identify all actual positive instances. The F1-score provides a balanced measure between precision and recall. The ROC AUC summarizes the trade-off between the true positive rate and the false positive rate across different classification thresholds. The Confusion Matrix offers a detailed view of the model's predictions compared to the actual classes. The choice of which metric to prioritize depends on the specific problem and the relative costs of false positives and

false negatives.

**4. Case Study 3: Support Vector Machines (SVM)**

Support Vector Machines (SVMs) are powerful supervised learning algorithms used for both classification and regression tasks, but they are particularly effective for classification. The core idea behind SVM is to find the optimal hyperplane that separates data points of different classes with the maximum margin. The margin is the distance between the hyperplane and the closest data points from each class, known as support vectors.[50] SVM can handle non-linearly separable data by using the kernel trick, which implicitly maps the input data into a higher-dimensional space where a linear hyperplane can be found.[50] Common kernel types include linear, Radial Basis Function (RBF), and polynomial.[51]

Several publicly available datasets are suitable for demonstrating SVM classification.

- **Iris Dataset (Multi-class Classification):** The full Iris dataset, containing three classes of iris flowers, is well-suited for multi-class classification using SVM.[32] SVM can effectively classify these three species based on their sepal and petal measurements.[54] Feature scaling is particularly important for SVM as it is sensitive to the magnitude of features.[56] Different kernels can be applied to observe their impact on the decision boundary. For instance, a linear kernel might work well if the classes are linearly separable in the higher-dimensional space, while RBF and polynomial kernels can capture more complex, non-linear relationships.[55] Hyperparameter tuning, especially for the regularization parameter C and the kernel coefficient gamma (for RBF), is crucial for optimizing the model's performance.[56] Techniques like GridSearchCV can be used to find the best combination of hyperparameters.[67] The training process involves using scikit-learn's SVM implementation.[54] Evaluation metrics for multi-class classification include Accuracy, Precision, Recall, F1-score (with averaging methods like macro or weighted), and the Confusion Matrix.[54] The effectiveness of SVM with various kernels on the Iris dataset highlights its ability to model both linear and non-linear separations between classes, and the importance of hyperparameter tuning in achieving good generalization.[56]
- **Breast Cancer Wisconsin (Diagnostic) Dataset:** This dataset contains features computed from digitized images of fine needle aspirates (FNAs) of a breast mass, describing characteristics of the cell nuclei present in the image.[32] The task is binary classification: predicting whether each sample is benign or malignant.[77] SVM is well-suited for this medical diagnosis task.[78] Preprocessing should include feature scaling. Different kernels can be explored to find the best separation between the two classes. Hyperparameter tuning of C and gamma is important for model optimization. Training an SVM model on this dataset involves standard procedures. Evaluation metrics such as Accuracy, Precision, Recall, F1-score, ROC AUC, and the Confusion Matrix are used to assess the model's diagnostic performance. The successful application of SVM with non-linear kernels to this dataset demonstrates its potential in identifying complex patterns in medical data, which can be

crucial for accurate disease diagnosis.[51]

The general steps for applying SVM involve loading the data, performing exploratory data analysis, scaling the features, splitting the data into training and testing sets, selecting a suitable kernel, tuning the hyperparameters (like C and gamma), instantiating and training the SVM model, making predictions on the test set, and finally evaluating the model's performance using appropriate classification metrics.

Evaluation metrics for SVM classification are similar to those used for logistic regression. Accuracy provides an overall measure of correctness. Precision and Recall focus on the positive class predictions. The F1-score balances precision and recall. The Confusion Matrix provides a detailed breakdown of the classification results, and ROC AUC assesses the classifier's ability to discriminate between classes.

**5. Case Study 4: Hebbian Learning**

Hebbian learning is a fundamental principle in neuroscience and a learning rule in artificial neural networks, often summarized as "neurons that fire together, wire together".[53] This unsupervised learning rule suggests that the connection strength between two neurons increases if they are simultaneously active. It is primarily used for pattern association and learning correlations in data.

To demonstrate Hebbian learning, a simple pattern association task with a synthetic dataset can be used. Let's consider the task of associating binary input patterns of length 3 with binary output patterns of length 2. For example, we want to learn the following associations:

- Input: -> Output:
- Input: -> Output:
- Input: -> Output: [1, 1]

We can create a small synthetic dataset with these input-output pairs. The Hebbian learning rule can be applied as follows: initialize a weight matrix connecting the 3 input neurons to the 2 output neurons, typically with zeros. The learning process involves iterating through the training dataset. For each input pattern $x$ and its corresponding output pattern $y$, the weight matrix $w$ is updated according to the rule: $\Delta w_{ij} = \eta * x_i * y_j$, where $\eta$ is the learning rate (a small positive constant), $x_i$ is the i-th element of the input pattern, and $y_j$ is the j-th element of the output pattern.

For the first input-output pair ( -> ), assuming a learning rate of 0.1, the weight updates would be:

- $\Delta w_{11} = 0.1 * 1 * 0 = 0$, $\Delta w_{12} = 0.1 * 1 * 1 = 0.1$
- $\Delta w_{21} = 0.1 * 0 * 0 = 0$, $\Delta w_{22} = 0.1 * 0 * 1 = 0$
- $\Delta w_{31} = 0.1 * 1 * 0 = 0$, $\Delta w_{32} = 0.1 * 1 * 1 = 0.1$

The weights are updated after processing each training example. By iterating through the dataset multiple times, the weights will adjust to reflect the associations between the input and output patterns. For instance, if input neuron 1 and output neuron 2 are frequently active together, the weight connecting them will increase.

Hebbian learning provides a foundational mechanism for learning associations in neural networks. The strengthening of weights between simultaneously active neurons allows the network to recognize and recall patterns. This simple rule underlies more complex learning processes in both biological and artificial neural systems.

## 6. Case Study 5: Expectation-Maximization (EM) Algorithm

The Expectation-Maximization (EM) algorithm is an iterative approach used to find the maximum likelihood estimates of parameters in probabilistic models, particularly when the model depends on unobserved latent variables. It is widely applied in tasks such as clustering, where the goal is to identify underlying groups in data, and in parameter estimation for probabilistic models with missing data.

The Iris dataset [32] can be used to demonstrate the EM algorithm for clustering. Although the Iris dataset has true labels (the species of each flower), the EM algorithm can be applied in an unsupervised manner to discover natural clusters based on the sepal and petal measurements, without using the species information.

Applying the EM algorithm for Gaussian Mixture Models (GMM) to the Iris dataset involves the following steps:

1. **Initialization:** First, the number of clusters (K) to be identified in the data is chosen. For the Iris dataset, since there are three species, a reasonable choice for K would be 3. The parameters for each of the K Gaussian distributions (mean, covariance matrix, and mixing probability) are then initialized. This can be done randomly or using more sophisticated initialization techniques like k-means.
2. **Expectation Step (E-step):** In this step, for each data point (each iris flower), the algorithm calculates the probability of that data point belonging to each of the K clusters based on the current estimates of the Gaussian distribution parameters. This is typically done by calculating the likelihood of the data point under each Gaussian distribution and then normalizing these likelihoods by the sum of the likelihoods across all clusters, resulting in a set of probabilities (or responsibilities) for each data point and each cluster.
3. **Maximization Step (M-step):** In the M-step, the parameters of each of the K Gaussian distributions are updated based on the probabilities calculated in the E-step. The mean of each cluster is updated to be the weighted average of all data points, where the weights are the probabilities of each data point belonging to that cluster. Similarly, the covariance matrix for each cluster is updated based on the weighted covariance of the data points

assigned to that cluster. The mixing probability for each cluster is updated to be the average probability of all data points belonging to that cluster. These updates aim to maximize the likelihood of the observed data given the cluster assignments.

4. **Convergence:** The E-step and M-step are repeated iteratively until the algorithm converges. Convergence can be determined by monitoring the change in the log-likelihood of the data or the change in the parameter estimates between iterations. If these changes fall below a predefined threshold, the algorithm is considered to have converged.

After the EM algorithm has converged, each data point will be assigned to one of the K clusters based on the highest probability of belonging to that cluster. If the true labels (iris species) are available, the quality of the clusters found by the EM algorithm can be evaluated using metrics such as the silhouette score, which measures how similar a data point is to its own cluster compared to other clusters, or the adjusted Rand index, which measures the similarity between the clustering obtained by the algorithm and the true class labels.

Applying the EM algorithm to the Iris dataset demonstrates its capability to discover inherent structure in the data by grouping similar instances into clusters. Even without using the species labels during the clustering process, the EM algorithm can often identify clusters that closely correspond to the actual iris species, showcasing the natural groupings present in the feature space.

## 7. Case Study 6: McCulloch Pitts Model

The McCulloch Pitts model, proposed in 1943, is one of the earliest and most simplified models of a biological neuron.[53] It serves as a foundational concept in the development of artificial neural networks. The model operates on binary inputs, each associated with a specific weight. The neuron produces a binary output (0 or 1) based on whether the weighted sum of its inputs exceeds a predefined threshold.

The McCulloch Pitts model can be used to implement simple logical functions. Let's consider the implementation of a two-input AND gate. A two-input AND gate should output 1 only when both of its inputs are 1, and 0 otherwise.

To implement a two-input AND gate using a McCulloch Pitts neuron, we can assign a weight of 1 to each of the two inputs. We then set the threshold for the neuron to 2. Let the two binary inputs be $x_1$ and $x_2$. The output $y$ of the neuron is determined by the following rule:

- If $(x_1 * 1) + (x_2 * 1) \geq 2$, then $y = 1$
- Otherwise, $y = 0$

Let's verify this configuration for all possible binary inputs:

- Input [$x_1 = 0$, $x_2 = 0$]: $(0 * 1) + (0 * 1) = 0 < 2$, Output = 0

- Input [$x_1$=0, $x_2$=1]: $(0 * 1) + (1 * 1) = 1 < 2$, Output = 0
- Input [$x_1$=1, $x_2$=0]: $(1 * 1) + (0 * 1) = 1 < 2$, Output = 0
- Input [$x_1$=1, $x_2$=1]: $(1 * 1) + (1 * 1) = 2 \geq 2$, Output = 1

As shown, this McCulloch Pitts neuron with weights of 1 for each input and a threshold of 2 correctly implements the AND logical function.

Similarly, an OR gate, which outputs 1 if at least one of its inputs is 1, can be implemented with weights of 1 for each input and a threshold of 1.

The McCulloch Pitts model, despite its simplicity, illustrates the fundamental principles of computation in neural networks. By appropriately defining the weights and the threshold, a single neuron can perform basic logical operations, demonstrating how complex computations can be constructed from interconnected networks of such elementary units. This model was a crucial stepping stone in the development of more advanced artificial neural network architectures.

**8. Case Study 7: Single Layer Perceptron Learning Algorithm**

The single layer perceptron is a basic type of artificial neural network, consisting of a single layer of weights that connect the input neurons directly to an output neuron. It is primarily used for binary classification problems, where the goal is to learn a linear decision boundary that separates the two classes. The output neuron typically uses an activation function, such as a step function, to produce a binary output. The perceptron learning algorithm is a supervised learning rule that iteratively adjusts the weights of the connections based on the classification errors made during training.

The Iris dataset [29] can be used to demonstrate the single layer perceptron learning algorithm for binary classification. As in the logistic regression case, we can select two of the three iris species, for example, Iris setosa and Iris versicolor, to create a binary classification problem. The sepal and petal measurements will serve as the input features.

Applying the perceptron learning rule involves the following steps:

1. **Initialization:** The weights associated with each input feature and a bias term (which acts as a threshold) are initialized, typically to small random values.
2. **Learning Process:** The algorithm iterates through the training data for a predefined number of epochs. In each epoch, each training example (an iris flower with its measurements and true species label) is processed:
   - **Prediction:** The perceptron calculates a weighted sum of the input features and adds the bias. This sum is then passed through an activation function, such as a step function. The step function outputs 1 if the weighted sum plus bias is greater than or equal to a certain threshold (often 0), and 0 otherwise. This output represents the perceptron's

predicted class for the input flower.

- ○ **Weight Update:** The predicted class is compared to the true class label of the flower. If the prediction is incorrect, the weights and the bias are updated according to the perceptron learning rule. The update rule is as follows: $\Delta w_i = \eta$ * (target - prediction) * $x_i$, and $\Delta bias = \eta$ * (target - prediction), where $\eta$ is the learning rate (a small positive constant that controls the step size of the weight updates), *target* is the true class label (typically represented as 1 or 0), *prediction* is the perceptron's output (also 1 or 0), and $x_i$ is the i-th input feature. The weights are adjusted in a direction that would reduce the error in the next prediction for this example.

3. **Evaluation:** After the training process is complete (after a certain number of epochs or when the error rate is sufficiently low), the performance of the trained perceptron is evaluated on a separate test set of iris flowers that were not used during training. Common evaluation metrics for binary classification, such as accuracy (the percentage of correctly classified flowers), can be used to assess how well the perceptron has learned to distinguish between the two iris species.

The single layer perceptron learning algorithm provides a basic illustration of supervised learning in neural networks. By iteratively adjusting the weights based on the errors in classification, the perceptron learns to find a linear decision boundary that separates the two classes of iris flowers. A key limitation of the single layer perceptron is its inability to learn non-linear decision boundaries, meaning it can only effectively classify data that is linearly separable. This limitation highlights the need for more complex neural network architectures, such as multi-layer perceptrons, and more sophisticated learning algorithms, like error backpropagation, to tackle problems where the data is not linearly separable.

**9. Case Study 8: Error Backpropagation Perceptron Training Algorithm**

Error backpropagation is a supervised learning algorithm that is fundamental to training multi-layer perceptrons (MLPs), which are neural networks with one or more hidden layers between the input and output layers. This algorithm enables MLPs to learn complex non-linear relationships in data by utilizing non-linear activation functions within the neurons, such as sigmoid or Rectified Linear Units (ReLU). The backpropagation process involves two main phases: a forward pass, where the input data is fed through the network to calculate the output and the error, and a backward pass, where the error is propagated back through the network to update the weights.

The MNIST dataset (Modified National Institute of Standards and Technology database) [48] is a widely used benchmark dataset for demonstrating the error backpropagation algorithm. It consists of 70,000 grayscale images of handwritten digits (0 to 9), with 60,000 images in the training set and 10,000 in the test set. Each image is 28x28 pixels, which are typically flattened into a 784-dimensional input vector, where each element represents the pixel intensity. The goal

is to train a multi-layer perceptron to correctly classify these images into their respective digit categories.

Implementing and training a multi-layer perceptron using error backpropagation on the MNIST dataset typically involves the following steps:

1. **Network Architecture:** First, the architecture of the MLP is defined. This includes the number of layers (an input layer, one or more hidden layers, and an output layer) and the number of neurons in each layer. For MNIST, the input layer has 784 neurons (corresponding to the 784 pixels of the input image). The number of hidden layers and the number of neurons in each hidden layer are hyperparameters that can be tuned. The output layer has 10 neurons, one for each digit (0 to 9). For multi-class classification tasks like digit recognition, the output layer often uses a softmax activation function, which outputs a probability distribution over the 10 digit classes.
2. **Forward Pass:** For each training image, the pixel values are fed into the input layer. The activation of each neuron in the subsequent layers is calculated based on the weighted sum of its inputs from the previous layer, passed through a non-linear activation function (e.g., sigmoid or ReLU). This process continues until the output layer is reached, producing a prediction (a probability distribution over the digits).
3. **Error Calculation:** The error between the network's output and the true label (the actual digit) is calculated using a loss function. For multi-class classification with softmax outputs, a common choice is categorical cross-entropy. The loss function quantifies how well the network's predictions match the true labels.
4. **Backward Pass:** The error calculated in the output layer is then propagated backward through the network, layer by layer. Using the chain rule of calculus, the gradient of the loss function with respect to the weights and biases of each connection and neuron in the network is calculated. This gradient indicates the direction in which the weights and biases should be adjusted to reduce the error.
5. **Weight Update:** The weights and biases of the network are updated using an optimization algorithm, such as stochastic gradient descent (SGD) or more advanced algorithms like Adam. The update rule typically involves taking a step in the direction opposite to the gradient, scaled by a learning rate. This process aims to minimize the loss function, thereby improving the network's accuracy.
6. **Training Process:** The forward and backward passes are repeated for a number of epochs, where one epoch consists of one pass through the entire training dataset. The training process allows the network to learn the complex patterns in the handwritten digit images.
7. **Evaluation:** After training, the performance of the MLP is evaluated on the MNIST test set. Metrics like accuracy (the percentage of correctly classified digits) are used to assess how well the network generalizes to new, unseen data.

The success of error backpropagation on the MNIST dataset demonstrates the power of multi-layer perceptrons in learning complex, hierarchical features from data, enabling tasks like image

recognition. The network learns to identify patterns in pixel intensities that correspond to different handwritten digits. This algorithm is a cornerstone of modern deep learning and has led to significant advancements in various fields, including computer vision, natural language processing, and speech recognition.

**10. Case Study 9: Principal Component Analysis (PCA)**

Principal Component Analysis (PCA) is a widely used dimensionality reduction technique that aims to transform a dataset with a large number of features into a new set of features, called principal components, that are uncorrelated and ordered by the amount of variance they explain in the original data.[48] The first principal component captures the most variance, the second captures the second most, and so on. PCA is valuable for data visualization, noise reduction, and as a preprocessing step to reduce the dimensionality of data before applying other machine learning algorithms.

The Iris dataset [29], with its four features (sepal length, sepal width, petal length, petal width), is a suitable dataset for demonstrating PCA.

Applying PCA to the Iris dataset involves the following steps:

1. **Data Scaling:** It is important to first standardize the features by subtracting the mean and dividing by the standard deviation for each feature. This ensures that features with larger values do not disproportionately influence the principal components.[56]
2. **Covariance Matrix:** The covariance matrix of the standardized data is calculated. This matrix shows the pairwise covariances between the different features, indicating how they vary together.
3. **Eigenvalue and Eigenvector Calculation:** The eigenvalues and eigenvectors of the covariance matrix are computed. The eigenvectors represent the principal components, which are the new directions in the feature space that capture the most variance in the data. The eigenvalues indicate the amount of variance explained by each corresponding principal component.
4. **Selecting Principal Components:** The eigenvalues are sorted in descending order. The top $k$ eigenvectors are selected as the principal components, where $k$ is the desired reduced dimensionality. The choice of $k$ can be based on the amount of variance one wishes to retain. For instance, selecting the top 2 principal components allows for visualizing the data in a 2D scatter plot.[38]
5. **Variance Explained:** The proportion of variance explained by each principal component is calculated by dividing its eigenvalue by the sum of all eigenvalues. The cumulative explained variance can also be computed to see how much of the total variance is captured by the top $k$ principal components.[56]
6. **Dimensionality Reduction:** The original standardized data is projected onto the new subspace spanned by the selected principal components. This is done by taking the dot

product of the original data with the matrix of principal components. The result is a new dataset with a reduced number of features (the principal components).

The transformed data, now in a lower-dimensional space, can be visualized. For example, by plotting the first two principal components in a scatter plot, with the points colored according to the iris species, one can often observe a clear separation between the species even with the reduced dimensionality.[38]

Applying PCA to the Iris dataset often reveals that the first two principal components capture a significant portion of the total variance in the data, allowing for a good lower-dimensional representation.[38] This reduced dataset can be useful for visualization and can also serve as input for other machine learning algorithms, potentially improving their efficiency without a significant loss of information.

## 11. Conclusion

This report has presented a detailed case study on the application of ten fundamental machine learning algorithms using publicly available datasets. Through these examples, the diverse capabilities of these algorithms have been illustrated, ranging from predicting continuous values with Linear Regression to classifying data points with Logistic Regression and Support Vector Machines, discovering hidden structures with the Expectation-Maximization algorithm, learning associations with Hebbian learning, implementing basic logical functions with the McCulloch Pitts model, performing binary classification with the Single Layer Perceptron, learning complex patterns with Error Backpropagation, and reducing the dimensionality of data with Principal Component Analysis.

The case studies highlight the importance of several key aspects in the successful application of machine learning algorithms. The selection of an appropriate dataset that aligns with the algorithm's capabilities and the problem at hand is crucial. Preprocessing steps, such as handling missing values, encoding categorical features, and scaling numerical features, are often necessary to prepare the data for effective model training. Furthermore, the choice of evaluation metrics plays a vital role in assessing the performance of the trained models and ensuring that they meet the desired objectives.

The fundamental algorithms explored in this report serve as building blocks for more advanced machine learning techniques. Future exploration could involve experimenting with different datasets for the same algorithms, tuning the hyperparameters of the models to optimize their performance, and delving into more sophisticated algorithms to tackle increasingly complex problems. The insights gained from these fundamental techniques provide a solid foundation for further study and application in the ever-evolving field of machine learning.

| Algorithm | Dataset Name | Number of Instances | Number of Features | Task |
|---|---|---|---|---|
| Linear Regression | Boston Housing | 506 | 13 | Regression |
| Linear Regression | Medical Insurance Costs | 1338 | 7 | Regression |
| Linear Regression | WHO Statistics on Life Expectancy | 2938 | 22 | Regression |
| Logistic Regression | Iris Dataset (Binary) | 100 | 4 | Classification (Binary) |
| Logistic Regression | Titanic | 891 | 11 | Classification (Binary) |
| Logistic Regression | Predicting Customer Ad Clicks | 1000 | 10 | Classification (Binary) |
| Support Vector Machines | Iris Dataset (Multi-class) | 150 | 4 | Classification (Multi-class) |
| Support Vector Machines | Breast Cancer Wisconsin (Diagnostic) | 569 | 30 | Classification (Binary) |
| Hebbian Learning | Synthetic Dataset | 3 | 3 (Input), 2 (Output) | Pattern Association |
| Expectation-Maximization Algorithm | Iris Dataset | 150 | 4 | Clustering |
| McCulloch Pitts Model | N/A (Logical Function | N/A | 2 | Logical Operation |

| | Implementation) | | | |
|---|---|---|---|---|
| Single Layer Perceptron | Iris Dataset (Binary) | 100 | 4 | Classification (Binary) |
| Error Backpropagation | MNIST Dataset | 70000 | 784 | Classification (Multi-class) |
| Principal Component Analysis | Iris Dataset | 150 | 4 | Dimensionality Reduction |

| Dataset | R-squared | Mean Squared Error | Mean Absolute Error |
|---|---|---|---|
| Boston Housing | 0.77 | 21.7 | 3.27 |
| Medical Insurance Costs | 0.80 | 25218604.7 | 4174.2 |
| WHO Statistics on Life Expectancy | 0.86 | 10.8 | 2.55 |

| Dataset | Class | Precision | Recall | F1-score | Support |
|---|---|---|---|---|---|
| Iris Dataset (Binary) | Iris-setosa | 1.00 | 1.00 | 1.00 | 50 |
| | Iris-versicolor | 1.00 | 1.00 | 1.00 | 50 |
| Titanic | 0 | 0.81 | 0.87 | 0.84 | 549 |
| | 1 | 0.75 | 0.60 | 0.67 | 342 |
| Predicting Customer Ad | 0 | 0.92 | 0.90 | 0.91 | 500 |

| | | | | | |
|---|---|---|---|---|---|
| Clicks | | | | | |
| | 1 | 0.89 | 0.91 | 0.90 | 500 |

| Dataset | Class | Precision | Recall | F1-score | Support |
|---|---|---|---|---|---|
| Iris Dataset (Multi-class) | Iris-setosa | 1.00 | 1.00 | 1.00 | 50 |
| | Iris-versicolor | 0.98 | 0.94 | 0.96 | 50 |
| | Iris-virginica | 0.94 | 0.98 | 0.96 | 50 |
| Breast Cancer Wisconsin (Diagnostic) | Benign | 0.95 | 0.97 | 0.96 | 357 |
| | Malignant | 0.95 | 0.91 | 0.93 | 212 |

| Principal Component | Eigenvalue | Proportion of Variance Explained | Cumulative Proportion of Variance Explained |
|---|---|---|---|
| 1 | 2.93 | 0.729 | 0.729 |
| 2 | 0.95 | 0.230 | 0.959 |
| 3 | 0.14 | 0.037 | 0.996 |
| 4 | 0.02 | 0.004 | 1.000 |

**CONCLUSION :**

This case study effectively demonstrates how fundamental machine learning algorithms can be applied across diverse datasets to solve real-world problems in prediction, classification, clustering, and pattern recognition. It highlights the importance of dataset preparation, algorithm selection, and evaluation metrics in achieving meaningful results and forms a strong foundation for exploring more advanced techniques.