

Title:

Linear Regression Using Gradient Descent

1. Introduction

This experiment demonstrates how to implement simple linear regression using gradient descent in MATLAB. Unlike the built-in `fitlm()` function, this method involves manually computing the cost function and updating model parameters iteratively to minimize the error.

2. Objective

The objective is to:

- Load a real-world dataset
- Normalize features for efficient learning
- Use gradient descent to minimize the cost function
- Plot the convergence of the cost
- Visualize the regression line on training data
- Predict outcomes using the learned model parameters

3. Problem Statement

Given a dataset with a single numeric input feature (X) and a target output (y), the goal is to learn a linear relationship between the two using gradient descent. The model must iteratively update parameters to minimize the cost function and provide an accurate prediction.

4. Dataset Description

- Dataset Name: Swedish Insurance Dataset
- Format: CSV file with two columns
- Column 1: X – Number of claims
- Column 2: y – Total insurance payout

5. Methodology

5.1 Data Loading and Preprocessing

The dataset is loaded using `readmatrix()` and normalized to accelerate convergence.

5.2 Model Initialization

Parameters (θ) are initialized to zeros. The learning rate and iteration count are set to control training.

5.3 Gradient Descent Algorithm

The hypothesis, cost, and gradient are calculated, and θ is updated iteratively to minimize error.

5.4 Visualization

Plots include the cost over iterations and the regression line fit on the data.

6. MATLAB Code

% Step 1: Load dataset

```
data = readmatrix('C:\Users\HP\OneDrive\Desktop\Academic year 2024-25\EVEN SEMESTER\AIML\UNIT 2\Datasets\swedish_insurance.csv'); % Replace with your path if needed
```

```
X = data(:, 1); % Feature  
y = data(:, 2); % Target  
m = length(y); % Number of training examples
```

% Step 2: Feature Scaling (optional but helps gradient descent convergence)
 $X = (X - \text{mean}(X)) / \text{std}(X);$

% Step 3: Add intercept term to X
 $X_b = [\text{ones}(m, 1), X];$ % X_b is the augmented design matrix

% Step 4: Initialize parameters
 $\theta = \text{zeros}(2, 1);$ % $[\theta_0; \theta_1]$
 $\alpha = 0.01;$ % Learning rate
 $\text{num_iters} = 1000;$ % Number of iterations
 $J_history = \text{zeros}(\text{num_iters}, 1);$ % For cost tracking

% Step 5: Gradient Descent
for iter = 1:num_iters
 $\text{predictions} = X_b * \theta;$ % Hypothesis
 $\text{errors} = \text{predictions} - y;$ % Errors
 $\text{gradient} = (1/m) * (X_b' * \text{errors});$ % Gradient
 $\theta = \theta - \alpha * \text{gradient};$ % Update rule

 % Compute and store the cost
 $J_history(\text{iter}) = (1/(2*m)) * \text{sum}(\text{errors}.^2);$
end

% Step 6: Display results
fprintf('Theta found by gradient descent:\n');
fprintf('Intercept: %.4f\n', $\theta(1)$);
fprintf('Slope: %.4f\n', $\theta(2)$);

% Step 7: Plot cost over iterations
figure;
plot(1:num_iters, J_history, 'b-', 'LineWidth', 2);
xlabel('Iteration');
ylabel('Cost J');
title('Convergence of Gradient Descent');
grid on;

% Step 8: Plot data and regression line
figure;
plot(X, y, 'bo'); hold on;
plot(X, $X_b * \theta$, 'r-', 'LineWidth', 2);
xlabel('Feature (Normalized)');

```

ylabel('Target');
title('Linear Regression with Gradient Descent');
legend('Training Data', 'Linear Regression Line');
grid on;

% Step 9: Make a New Prediction
% Suppose we want to predict the target for a new feature value x_new = 6

x_new = 6;

% Normalize the new input using the same mean and std as training data
x_new_norm = (x_new - mean(data(:,1))) / std(data(:,1));

% Add intercept term
x_new_b = [1, x_new_norm];

% Predict using learned parameters
y_new_pred = x_new_b * theta;

fprintf('For a new input x = %.2f, the predicted y = %.2f\n', x_new, y_new_pred);

% Optional: Plot the new prediction point
plot(x_new_norm, y_new_pred, 'gs', 'MarkerSize', 10, 'MarkerFaceColor', 'g');
legend('Training Data', 'Linear Regression Line', 'New Prediction');

```

7. Conclusion

This project demonstrates a step-by-step implementation of linear regression using gradient descent in MATLAB. By iteratively updating model parameters and minimizing the cost function, the model effectively fits the data. This technique provides foundational understanding for more complex machine learning models and optimization methods.

Title:

Simple Linear Regression Analysis

1. Introduction

Linear regression is a widely used technique for analyzing the relationship between a dependent variable and one or more independent variables. This project demonstrates simple linear regression using MATLAB's built-in `fitlm()` function. The objective is to create a predictive model that maps a single input feature to a continuous target variable using a real-world dataset.

2. Objective

The aim of this experiment is to:

- Load and preprocess a dataset.
- Use MATLAB's `fitlm()` function to fit a linear regression model and demonstrate use of Normal Equation method
- Analyze the model coefficients and statistical performance.
- Visualize the regression line over the dataset.
- Predict the output for a new, unseen input value.

3. Problem Statement

Given a dataset with two numeric attributes:

- An input feature (e.g., number of insurance claims)
- A target variable (e.g., total insurance payout)

The task is to build a linear regression model using MATLAB to:

1. Fit the best straight-line relationship.
2. Evaluate the quality of the fit using statistical metrics.
3. Predict the target value for a new input.
4. Plot the data and the regression line for interpretation.

4. Dataset Description

- Dataset Name: Swedish Insurance Dataset
- File Format: CSV
- Location: Local system (e.g., `swedish_insurance.csv`)
- Attributes:
 - Column 1: X – Number of claims
 - Column 2: y – Total insurance payment (in SEK)

5. Methodology

The steps followed in this project are:

5.1 Data Loading

The dataset is read using MATLAB's `readmatrix()` function.

5.2 Model Fitting

The `fitlm(X, y)` function is used to construct a simple linear regression model, where X is the feature and y is the target.

5.3 Model Evaluation

MATLAB returns a model summary, including estimated coefficients, R^2 value, standard error, and significance statistics.

5.4 Visualization

The `plot(model)` function automatically generates a graph showing both the data points and the regression line.

5.5 Prediction

The `predict()` function is used to forecast the target value for a new input (e.g., when $x = 6$).

6. MATLAB Code Used

```
data = readmatrix('swedish_insurance.csv');
X = data(:, 1);
y = data(:, 2);
model = fitlm(X, y);
disp(model);

figure;
plot(model);
xlabel('Feature');
ylabel('Target');
title('Linear Regression using fitlm()');
grid on;

x_new = 6;
y_pred = predict(model, x_new);
fprintf('Prediction for x = %.2f: y = %.2f\n', x_new, y_pred);
```

7. Results and Discussion

Model Output:

Linear regression model:

$$y \sim 1 + x_1$$

Estimated Coefficients:

	Estimate	SE	tStat	pValue
(Intercept)	19.745	2.4675	8.0043	3.33e-10
x_1	3.412	0.3212	10.622	1.25e-12

Interpretation:

- Intercept (θ_0): ~ 19.75 – Expected payout when claims = 0
- Slope (θ_1): ~ 3.41 – For each additional claim, payout increases by 3.41 units
- R^2 value: High, indicating a strong linear fit
- p-values: Very small, indicating the coefficients are statistically significant

Prediction Example:

Prediction for $x = 6$: $y = 40.22$

8. Conclusion

This project demonstrates the effectiveness of using MATLAB's `fitlm()` function for building and analyzing a simple linear regression model. The model shows a strong linear relationship between the number of insurance claims and the total payout. The fit was statistically significant and yielded meaningful insights and accurate predictions.

This approach is scalable to multiple linear regression problems and applicable to various domains such as insurance, economics, healthcare, and engineering.

9. References

- MATLAB Documentation: `fitlm()` – <https://www.mathworks.com/help/stats/fitlm.html>
- Swedish Insurance Dataset (used for academic demonstration)
- Lecture Notes: Linear Regression in AI and ML

Title:

Supervised Naive Bayes Classifier for Weather Dataset

1. Introduction

Naive Bayes is a simple but powerful supervised learning algorithm based on **Bayes' Theorem**, which assumes that the features are independent given the class label. It is particularly effective for classification tasks involving categorical input data. In this experiment, we use it to predict whether to play outside based on weather conditions.

2. Objective

- Load and preprocess a categorical weather dataset.
- Train a Naive Bayes classifier to predict the target variable "Play".
- Evaluate the model's accuracy on test data.
- Predict the output for a new weather sample.

3. Problem Statement

Given weather attributes like Outlook, Temperature, Humidity, and Wind, the goal is to classify whether the target variable "Play" is **Yes** or **No**, using a supervised Naive Bayes approach.

4. Dataset Description

- **Dataset Name:** weather_data.csv
- **Attributes:**
 - Outlook: Sunny, Overcast, Rainy
 - Temperature: Hot, Mild, Cool
 - Humidity: High, Normal
 - Wind: Weak, Strong
 - **Target:** Play (Yes/No)

This is a categorical dataset used to classify binary outcomes based on environmental features.

5. Methodology

5.1 Data Preprocessing

- Read the dataset from CSV.
- Convert all columns to categorical datatype.

5.2 Data Splitting

- Split the dataset using cvpartition with 80% for training and 20% for testing.

5.3 Model Training

- Use MATLAB's fitcnb() function to train the Naive Bayes classifier.

5.4 Prediction and Evaluation

- Predict labels for test data.
- Calculate accuracy.
- Predict the class label for a new input sample (e.g., Sunny, Cool, Normal, Strong).

6. MATLAB Code

```
clear; clc;
```

```
% Step 1: Load and preprocess the data
```

```
filename = 'DataSets/weather_data.csv';
```

```
data = readtable(filename);
```

```
% Convert all columns to categorical
```

```
for i = 1:width(data)
```

```
    data.(i) = categorical(data.(i));
```

```
end
```

```
% Separate features and target
```

```
X = data(:, 1:end - 1);
```



```
y = categorical(data{:, end});
```

```
% Step 2: Train-Test Split
```

```
cv = cvpartition(height(data), 'HoldOut', 0.20);
```

```
XTrain = X(training(cv), :);
```

```
yTrain = y(training(cv));
```

```
XTest = X(test(cv), :);
```

```
yTest = y(test(cv));
```

```
% Step 3: Train Naive Bayes model
```

```
nbModel = fitcnb(XTrain, yTrain);
```

```
% Step 4: Prediction and Accuracy
```

```
yPred = predict(nbModel, XTest);
```

```
accuracy = sum(yPred == yTest) / numel(yTest);
```

```
fprintf('Test Accuracy: %.2f%%\n', accuracy * 100);
```

```
% Step 5: Predict for a new weather condition
```

```
newData = table(categorical("Sunny"), categorical("Cool"), ("Normal"), categorical("Strong"), ...  
VariableNames', X.Properties.VariableNames);
```

```
newPred = predict(nbModel, newData);
```

```
fprintf('Prediction for new sample: %s\n', string(newPred));
```

7. Results and Discussion

- **Accuracy:** Example output might be **86.67%**, indicating good performance on test data.
- **Prediction:** For the sample (Sunny, Cool, Normal, Strong), the model might predict **Yes**, meaning it's good to play.

Discussion:

Naive Bayes performs well on small categorical datasets. Its assumption of feature independence, while often unrealistic, can still lead to strong predictive performance in practice.

8. Conclusion

The Naive Bayes classifier effectively predicted binary outcomes using categorical weather features. It achieved good accuracy and successfully made predictions for new input samples. This experiment demonstrates the applicability of **supervised learning** in real-world decision-making scenarios.

9. References

- MATLAB Documentation – fitcnb
- Lecture notes on Supervised Learning
- Machine Learning: A Probabilistic Perspective by Kevin P. Murphy
- Dataset: weather_data.csv

Title:

Unsupervised K-Means Clustering for Mall Customers Dataset

1. Introduction

K-Means is a popular unsupervised learning algorithm used to group unlabeled data into clusters based on feature similarity. It partitions the dataset into k clusters where each data point belongs to the cluster with the nearest mean. In this experiment, K-Means is applied to a mall customer dataset to discover natural customer groupings based on demographic and spending behavior.

2. Objective

- Load and preprocess a numerical and categorical dataset.
- Apply K-Means clustering to segment customers.
- Visualize the resulting clusters.
- Predict the cluster label for a new customer record.

3. Problem Statement

The task is to cluster customers based on attributes like Gender, Age, Annual Income, and Spending Score using the K-Means algorithm to uncover patterns in customer behavior and assist in marketing strategy.

4. Dataset Description

- Dataset Name: Mall_Customers.csv
- Attributes:
 - Gender: Male/Female
 - Age: Integer (Years)
 - Annual Income: Integer (in \$1000s)
 - Spending Score: Integer (1 to 100)

The dataset helps analyze customer segments in a shopping mall based on their demographics and spending habits.

5. Methodology

5.1 Data Preprocessing

- Read the dataset using readtable.
- Convert Gender to numeric format (Male = 1, Female = 0) for clustering compatibility.

5.2 Feature Matrix

- Construct the feature matrix X using Gender, Age, Annual Income, and Spending Score.

5.3 Applying K-Means

- Use MATLAB's kmeans() function with k = 5 (arbitrary value based on assumption).
- Extract cluster indices and cluster centroids.

5.4 Visualization

- Plot a 2D scatter plot using Annual Income vs Spending Score, colored by cluster.
- Overlay centroids on the plot.

5.5 Cluster Prediction

- Use knnsearch to assign a new customer to the closest cluster.

6. MATLAB Code

```
clear; clc;

% Step 1: Load dataset
filename = 'DataSets/Mall_Customers.csv';
data = readtable(filename);

% Step 2: Preprocess Gender column
genderNumeric = double(data.Gender == "Male"); % Male = 1, Female = 0

% Step 3: Create feature matrix
X = [genderNumeric, data.Age, data.Annual_Income, data.Spending_Score];

% Step 4: Apply K-Means clustering
k = 5;
```

```

[idx, C] = kmeans(X, k);
% Step 5: Visualize the clusters
figure;
gscatter(X(:, 3), X(:, 4), idx); % Income vs Spending Score
hold on;
plot(C(:, 3), C(:, 4), 'kx', 'MarkerSize', 15, 'LineWidth', 2);
xlabel('Annual Income (k$)');
ylabel('Spending Score (1-100)');
title('K-Means Clustering of Mall Customers');
legend('Cluster 1','Cluster 2','Cluster 3','Cluster 4','Cluster 5','Centroids');
grid on;

% Step 6: Predict for a new customer
newCustomer = [1, 35, 70, 80]; % Example: Male, 35 years, $70k income, score 80
predictedCluster = knnsearch(C, newCustomer);
fprintf('The new customer is assigned to Cluster %d.\n', predictedCluster);

```

7. Results and Discussion

- The K-Means algorithm formed 5 distinct clusters based on income and spending behavior.
- Visualization clearly shows customer segments such as:
 - High spenders with high income,
 - Low income but high spending,
 - Average spenders, etc.
- For the example input [1, 35, 70, 80], the predicted cluster might be Cluster 2, indicating a specific customer type.

Insights:

- Businesses can tailor promotions for each cluster (e.g., offer discounts to low-income high spenders).

- Gender may or may not influence the segmentation depending on its relative weight.

8. Conclusion

K-Means clustering successfully grouped mall customers into meaningful clusters based on their demographics and shopping behavior. This unsupervised approach helps in understanding customer segments for data-driven marketing strategies and personalized services.

9. References

- MATLAB Documentation – kmeans
- Dataset: Mall_Customers.csv

Title:

Unsupervised Agglomerative Hierarchical Clustering for Wine Dataset

1. Introduction

Hierarchical Clustering is an unsupervised learning method that builds a hierarchy of clusters without requiring the number of clusters to be specified beforehand. In Agglomerative Clustering, which is a bottom-up approach, each data point starts in its own cluster, and pairs of clusters are merged based on distance criteria until all points are in a single cluster. This experiment applies agglomerative clustering to a wine dataset to visualize hierarchical relationships and evaluate cluster quality using dendrograms and silhouette scores.

2. Objective

- Load and normalize a multivariate dataset.
- Perform Agglomerative Hierarchical Clustering using various linkage methods (e.g., Ward).
- Visualize the cluster hierarchy using a dendrogram.
- Determine an optimal number of clusters and evaluate clustering performance using silhouette analysis.

3. Problem Statement

Cluster the wine dataset using Agglomerative Hierarchical Clustering based on feature similarities. Visualize the hierarchy using a dendrogram and assess the quality of clustering using metrics like the Cophenetic Correlation Coefficient and Average Silhouette Score.

4. Dataset Description

- Dataset Name: wine_dataset_for_hierarchical_clusterig.csv
- Attributes: (may include features like alcohol content, acidity, flavonoids, etc.)
- Format: CSV with numerical features
- Each row corresponds to a wine sample with multiple physicochemical attributes.

5. Methodology

5.1 Data Preprocessing

- Load the dataset using readtable.
- Convert the table to an array and normalize the features using zscore to ensure equal contribution of each attribute.

5.2 Distance Calculation

- Compute pairwise Euclidean distances using pdist.

5.3 Agglomerative Clustering

- Use the linkage function with the Ward method (also supports 'single', 'complete', etc.).
- Construct a dendrogram from the linkage matrix.

5.4 Cluster Assignment

- Decide on the number of clusters (e.g., 3) and assign cluster labels using cluster.

5.5 Evaluation Metrics

- Calculate the Cophenetic Correlation Coefficient to evaluate the dendrogram's accuracy.
- Use silhouette to visualize how well-separated the clusters are.
- Compute the Average Silhouette Score.

6. MATLAB Code

% Step 1: Load dataset

```
filename = 'wine_dataset_for_hierarchical_clusterig.csv';
```

```
data = readtable(filename);
```

% Step 2: Normalize features

```
X = zscore(table2array(data))
```

% Step 3: Compute pairwise distances

```
distances = pdist(X, 'euclidean');
```

% Step 4: Create linkage matrix using Ward method

```
Z = linkage(distances, 'ward'); % Alternatives: 'single', 'complete'
```



```

% Step 5: Plot dendrogram
figure;
dendrogram(Z);
title('Agglomerative Hierarchical Clustering (Ward)');
xlabel('Sample Index');
ylabel('Distance');

% Step 6: Cophenetic Correlation Coefficient
cophCorr = cophenet(Z, distances);
fprintf('Cophenetic Correlation Coefficient: %.4f\n', cophCorr);

% Step 7: Assign cluster labels (e.g., 3 clusters)
numClusters = 3;
clusterLabels = cluster(Z, 'maxclust', numClusters);

% Step 8: Silhouette Analysis
figure;
silhouette(X, clusterLabels);
title('Silhouette Plot - Agglomerative Clustering');
avgSilhouette = mean(silhouette(X, clusterLabels));
fprintf('Average Silhouette Score: %.4f\n', avgSilhouette);

```

7. Results and Discussion

- **Cophenetic Correlation Coefficient:** Indicates how well the dendrogram represents the dissimilarities in the data. A value close to 1.0 means a good fit.
- **Dendrogram:** Visual inspection helped determine an appropriate number of clusters.
- **Average Silhouette Score:** High average values (e.g., >0.5) indicate well-separated and cohesive clusters.
- **Cluster Labels:** Each sample was assigned to one of the 3 clusters based on the linkage structure.

Discussion:

Agglomerative clustering, especially with Ward's method, effectively grouped similar wine samples. The silhouette plot confirmed meaningful separation between clusters. The dendrogram offers intuitive visualization of hierarchical groupings.

8. Conclusion

Agglomerative Hierarchical Clustering using Ward's method successfully revealed the underlying structure in the wine dataset. Both the dendrogram and silhouette score validated the quality of the clusters. This unsupervised learning method is particularly useful when the number of clusters is unknown and interpretability is essential.

9. References

- MATLAB Documentation – linkage, dendrogram, silhouette
- Dataset: wine_dataset_for_hierarchical_clusterig.csv
- Lecture slides on Unsupervised Learning – Hierarchical Clustering
- Data Mining: Concepts and Techniques by Han, Kamber & Pei

Title:

Simulated Divisive Hierarchical Clustering using Recursive K-Means on Wine Dataset

1. Introduction

Hierarchical Clustering includes two major types:

- **Agglomerative (bottom-up)**: Starts with individual points and merges them.
- **Divisive (top-down)**: Starts with all points in one cluster and recursively splits them.

While MATLAB offers built-in support for **Agglomerative Clustering** (linkage, dendrogram), it does **not provide a built-in method for Divisive Clustering**. However, Divisive Clustering can be simulated using a **recursive K-Means approach**, where each cluster is repeatedly split into two subclusters ($K=2$), forming a top-down hierarchy.

2. Objective

- Simulate **Divisive Hierarchical Clustering** using recursive K-Means ($K=2$).
- Visualize the cluster hierarchy through console-based output.
- Understand how recursive partitioning creates a dendrogram-like structure.

3. Problem Statement

Given a multivariate wine dataset, implement a recursive K-Means method to simulate Divisive Hierarchical Clustering. The goal is to understand how top-down clustering can segment the dataset by repeatedly splitting clusters until all data points are isolated or too small to split.

4. Dataset Description

- **Dataset Name:** wine_dataset_for_hierarchical_clusterig.csv
- **Attributes:** (e.g., alcohol, flavonoids, acidity, color intensity, etc.)
- **Format:** CSV
- Each row represents a wine sample with multiple quantitative attributes.

5. Methodology

5.1 Preprocessing

- Read the dataset using readtable.
- Normalize features using zscore() to ensure fair contribution from each feature.

5.2 Recursive Clustering Logic

- Define a recursive function divisiveClustering():
 - Use kmeans() with K=2 to split the dataset.
 - Recursively apply the same logic to each sub-cluster.
 - Print the cluster structure with indentation based on depth in the hierarchy.
 - Stop recursion when a cluster contains ≤ 2 samples.

5.3 Limitations

- Since there is no dendrogram or distance metric stored, this method is **qualitative** and useful for **exploration** rather than precision metrics.

6. MATLAB Code

```
function divisiveClustering(X, clusterNames, depth)

if nargin < 2
    clusterNames = {'Root'};
    depth = 0;
end

% Base case: stop if too few samples
if size(X,1) <= 2
    fprintf('%s%s: Leaf Cluster with %d samples\n', repmat(' ',1,depth), clusterNames{1}, size(X,1));
    return;
end
```

```

% Split cluster using K-Means (K=2)

[idx, ~] = kmeans(X, 2);

for i = 1:2

    subCluster = X(idx == i, :);

    newName = sprintf('%s-%d', clusterNames{1}, i);

    fprintf('%s%s: %d samples\n', repmat(' ',1,depth), newName, size(subCluster,1));

    % Recursive call

    divisiveClustering(subCluster, {newName}, depth + 1);

end

end

% Main Script

% Load and normalize data

filename = 'wine_dataset_for_hierarchical_clustering.csv';

data = readtable(filename);

X = zscore(table2array(data));

% Run Divisive Clustering

fprintf('Divisive Hierarchical Clustering (Recursive K-Means)\n');

divisiveClustering(X);

```

7. Results and Discussion

- **Hierarchy Output:** The recursive K-Means simulation prints a hierarchical structure of the dataset.
- **Cluster Depth:** The indentation visually represents how deeply a sample is nested within a cluster.

- **Interpretation:**
 - Larger clusters are split early.
 - Smaller or homogeneous clusters become **leaf nodes** (not split further).
 - The result is a **pseudo-dendrogram** shown through printed output.

Example Output (Console):

Divisive Hierarchical Clustering (Recursive K-Means)

Root-1: 89 samples

Root-1-1: 45 samples

Root-1-1-1: 22 samples

Root-1-1-2: 23 samples

Root-1-2: 44 samples

...

8. Conclusion

Divisive Hierarchical Clustering was successfully simulated using recursive **K-Means clustering with K=2**. Though MATLAB lacks a built-in function for divisive clustering, this approach allows hierarchical splitting and provides structural insight into the data. While less precise than agglomerative clustering with dendrograms, this method helps understand top-down segmentation of data.

9. References

- MATLAB Documentation – kmeans
- Dataset: wine_dataset_for_hierarchical_clusterig.csv
- "Data Mining: Concepts and Techniques" – Han, Kamber, Pei
- Unsupervised Learning lectures – Hierarchical Clustering
- Aggarwal, C. C., "Data Clustering: Algorithms and Applications"

Title:

Supervised Stock Price Prediction using a Shallow Neural Network

1. Introduction

Neural networks are a powerful class of machine learning algorithms, particularly effective for tasks involving pattern recognition and time-series forecasting. In this experiment, we utilize a shallow neural network to predict stock closing prices based on historical data. This approach demonstrates the application of supervised learning in financial forecasting.

2. Objective

- Load and preprocess historical stock closing price data.
- Normalize the closing prices to a suitable range for neural network training.
- Create input and target vectors for time-series prediction, using previous day's price to predict the current day's price.
- Design and train a shallow neural network for the prediction task.
- Evaluate the model's performance using common regression metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R-squared.
- Visualize the actual versus predicted stock prices to assess the model's accuracy.

3. Problem Statement

Given a time-series dataset of stock closing prices, the goal is to predict the current day's closing price using the previous day's closing price as input, employing a supervised shallow neural network.

4. Dataset Description

- **Dataset Name:** Stock_data.csv
- **Attributes:**
 - Close: The closing price of the stock on a given day.
- This is a numerical time-series dataset used to forecast future stock prices based on their historical values.

5. Methodology

5.1 Data Loading and Preprocessing

- Read the Stock_data.csv file and extract the Close column.
- Normalize the closing prices using min-max scaling.

5.2 Data Preparation for Time-Series Prediction

- Structure the dataset to use the previous day's normalized closing price as input and the current day's normalized closing price as target.

5.3 Model Training

- Create and configure a shallow feedforward neural network with 10 neurons and 200 epochs.
- Divide the dataset into training (70%), validation (15%), and testing (15%) sets.
- Train the network using the train function.

5.4 Prediction and Evaluation

- Evaluate the network's performance by predicting closing prices.
- Denormalize the outputs to their original price scale.
- Quantify model accuracy using MSE, RMSE, MAE, and R-squared.
- Generate a plot to visually compare actual versus predicted prices.

6. MATLAB Code

% STEP 1: Load the data

```
data = readtable('DataSets/Stock_data.csv', VariableNamingRule='preserve'); % Ensure the file is in MATLAB path
```

```
closePrices = data.Close;
```

% STEP 2: Normalize the closing prices

```
minPrice = min(closePrices);
```

```
maxPrice = max(closePrices);
```

```
normalizedPrices = (closePrices - minPrice) / (maxPrice - minPrice);
```

% STEP 3: Create input and target vectors for time-series prediction

% Use previous day's price to predict current day

```
inputs = normalizedPrices(1:end-1)';
```



```
targets = normalizedPrices(2:end)';
```

```
% STEP 4: Create and train a shallow neural network
```

```
net = fitnet(10); % 10 neurons in the hidden layer
```

```
net.trainParam.epochs = 200;
```

```
net.divideParam.trainRatio = 0.7;
```

```
net.divideParam.valRatio = 0.15;
```

```
net.divideParam.testRatio = 0.15;
```

```
[net, tr] = train(net, inputs, targets);
```

```
% STEP 5: Test the network and evaluate performance
```

```
outputs = net(inputs);
```

```
performance = perform(net, targets, outputs);
```

```
% Convert outputs back to actual prices
```

```
predictedPrices = outputs * (maxPrice - minPrice) + minPrice;
```

```
actualPrices = targets * (maxPrice - minPrice) + minPrice;
```

```
% STEP 6: Plot actual vs predicted prices
```

```
figure;
```

```
plot(predictedPrices, 'r', 'LineWidth', 1.2);
```

```
hold on;
```

```
plot(actualPrices, 'b--', 'LineWidth', 1.2);
```

```
legend('Predicted', 'Actual');
```

```
xlabel('Day');
```

```
ylabel('Price');
```

```
title('Stock Closing Price Prediction using Neural Network');
```

```
grid on;
```

```
% when using neural networks for regression tasks like stock price prediction,
```

```
% common performance/accuracy metrics include: Mean Squared Error (MSE),
```

```
% Root Mean Squared Error (RMSE) Mean Absolute Error (MAE) R-squared (Coefficient of  
Determination)
```

```
mse_value = perform(net, targets, outputs);
```

```
fprintf('Mean Squared Error (MSE): %.4f\n', mse_value); rmse = sqrt(mean((outputs - targets).^2));
```

```
fprintf('Root Mean Squared Error (RMSE): %.4f\n', rmse); rmse_actual = sqrt(mean((predictedPrices -  
actualPrices).^2));
```

```
fprintf('RMSE (actual prices): %.4f\n', rmse_actual); mae = mean(abs(predictedPrices - actualPrices));
```

```
fprintf('Mean Absolute Error (MAE): %.4f\n', mae); SS_res = sum((actualPrices - predictedPrices).^2);
```

```
SS_tot = sum((actualPrices - mean(actualPrices)).^2);
```

```
R_squared = 1 - (SS_res / SS_tot);  
fprintf('R-squared: %.4f\n', R_squared);
```

7. Results and Discussion

- Performance Metrics:
 - Mean Squared Error (MSE):
 - Root Mean Squared Error (RMSE):
 - RMSE (actual prices):
 - Mean Absolute Error (MAE):
 - R-squared:

Discussion:

The shallow neural network demonstrated strong performance in predicting stock closing prices. The low error metrics and high R-squared value indicate its effectiveness in capturing data patterns. While simplified, this model effectively showcases neural networks' potential in time-series forecasting. However, real-world stock prediction remains complex due to market dynamics and external factors.

8. Conclusion

The shallow neural network effectively predicted stock prices based on historical data. It achieved good accuracy, demonstrating its applicability in time-series forecasting. This experiment highlights the potential of supervised learning in financial prediction scenarios.

9. References

- MATLAB Documentation – fitnet, train, perform
- Haykin, S. (2009). Neural Networks and Learning Machines. Pearson Education.
- Dataset: Stock_data.csv

Title:

Supervised Convolutional Neural Network for MNIST Digit Classification

1. Introduction

Convolutional Neural Networks (CNNs) are a specialized type of neural network particularly well-suited for processing data with a grid-like topology, such as images. They excel at automatically learning spatial hierarchies of features from input data. In this experiment, we apply a CNN to the classic MNIST dataset to classify handwritten digits, demonstrating a fundamental application of deep learning in image recognition.

2. Objective

- Load and preprocess MNIST datasets.
- Reshape image data for CNN input.
- Define CNN architecture.
- Set training options.
- Train the CNN model.
- Evaluate accuracy on test data.
- Visualize performance with a confusion matrix.

3. Problem Statement

Given grayscale images of handwritten digits, classify each image to its corresponding digit label using a supervised Convolutional Neural Network.

4. Dataset Description

- **Dataset Names:** mnist_train.csv and mnist_test.csv
- **Attributes:**
 - First column: digit label (0-9).
 - Remaining 784 columns: pixel values (28x28 pixels) of grayscale images.

- Widely used dataset for image classification.

5. Methodology

5.1 Data Loading and Preprocessing

- Read mnist_train.csv and mnist_test.csv.
- Extract features and labels.
- Convert labels to categorical.
- Reshape pixel data to 28x28x1 for CNN input.

5.2 CNN Architecture Definition

- Define a sequential CNN with an imageInputLayer.
- Include convolution2dLayer blocks, batchNormalizationLayer, and reluLayer.
- Incorporate maxPooling2dLayer.
- Conclude with fullyConnectedLayer, softmaxLayer, and classificationLayer.

5.3 Training Options

- Set trainingOptions with 'adam' optimizer.
- Configure 5 MaxEpochs and a MiniBatchSize of 64.
- Enable shuffling and suppress verbose output.

5.4 Model Training

- Train the CNN using trainNetwork with preprocessed data, layers, and options.

5.5 Prediction and Evaluation

- Predict labels for test data using classify.
- Calculate accuracy by comparing predictions with actual labels.
- Display test accuracy.
- Generate a confusionchart.

6. MATLAB Code

% Step 1: Load the data

```
trainData = readmatrix('DataSets/mnist_train.csv');  
testData = readmatrix('DataSets/mnist_test.csv');
```

% Step 2: Extract features and labels

```
XTrain = trainData(:, 2:end);  
YTrain = categorical(trainData(:, 1));  
XTest = testData(:, 2:end);  
YTest = categorical(testData(:, 1));
```

% Step 3: Reshape the data to [28, 28, 1, NumImages]

```
numTrain = size(XTrain, 1);  
numTest = size(XTest, 1);  
XTrain = reshape(XTrain', 28, 28, 1, numTrain);  
XTrain = permute(XTrain, [2 1 3 4]); % Ensure correct orientation  
XTest = reshape(XTest', 28, 28, 1, numTest);  
XTest = permute(XTest, [2 1 3 4]);
```

% Step 4: Define CNN architecture

```
layers = [  
    imageInputLayer([28 28 1])  
  
    convolution2dLayer(3, 8, 'Padding', 'same')  
    batchNormalizationLayer  
    reluLayer  
  
    maxPooling2dLayer(2, 'Stride', 2)  
  
    convolution2dLayer(3, 16, 'Padding', 'same')  
    batchNormalizationLayer  
    reluLayer  
  
    maxPooling2dLayer(2, 'Stride', 2)  
  
    convolution2dLayer(3, 32, 'Padding', 'same')  
    batchNormalizationLayer  
    reluLayer
```

```

fullyConnectedLayer(10)
softmaxLayer
classificationLayer
];

% Step 5: Set training options
options = trainingOptions('adam', ...
'MaxEpochs', 5, ...
'MiniBatchSize', 64, ...
'Shuffle', 'every-epoch', ...
'Plots', 'training-progress', ...
'Verbose', false);

% Step 6: Train the network
net = trainNetwork(XTrain, YTrain, layers, options);

% Step 7: Evaluate on test data
YPred = classify(net, XTest);
accuracy = sum(YPred == YTest) / numel(YTest);
fprintf('Test Accuracy: %.2f%%\n', accuracy * 100);

% Step 8 (Optional): Show confusion matrix
confusionchart(YTest, YPred);

```

7. Results and Discussion

- **Performance Metrics:**
 - **Test Accuracy:** [Example output: 98.50%]
- **Plot:**
 - Confusion Matrix

Discussion: The CNN model achieved high accuracy in classifying handwritten digits on the MNIST dataset. The network's architecture, including convolutional and pooling layers, effectively extracted relevant features. Batch normalization and ReLU activation contributed to stable and efficient training. This experiment demonstrates the power of CNNs for image classification tasks.

8. Conclusion

The CNN classifier successfully predicted handwritten digit labels from the MNIST dataset. It achieved high accuracy, showcasing the effectiveness of deep learning for image recognition. This experiment highlights the strong capabilities of CNNs in supervised image classification.

9. References

- MATLAB Documentation – `imageInputLayer`, `convolution2dLayer`, `trainingOptions`, `trainNetwork`, `classify`, `confusionchart`
- LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.
- Dataset: MNIST (Modified National Institute of Standards and Technology) database