

DL Lab Assignment-2

Name: Kanishk Vikram Singh

SAP ID: 500121743

Date: August 16, 2025

1. Potential Applications of Supervised Learning

Supervised learning is a subfield of machine learning where models are trained on labeled data. This means that for each data input, there is a corresponding correct output. The model learns the mapping function from the input to the output and can then make predictions on new, unseen data. Below are five potential applications.

a. Email Spam Detection (Classification)

This is one of the most common applications. Models are trained on a large dataset of emails that have been pre-labeled as either 'spam' or 'not spam' (ham).

How it's used: An email provider (like Gmail) uses this model to automatically filter incoming emails.

Input Features: Words in the email body and subject, sender's address, email metadata.

Output: A binary label: spam (1) or ham (0).

b. House Price Prediction (Regression)

This model predicts the continuous value of a house's selling price based on its features.

How it's used: Real estate companies (like Zillow or Redfin) use these models to provide users with an estimated market value for properties.

Input Features: Square footage, number of bedrooms/bathrooms, location, age of the house, presence of a garage.

Output: A continuous value representing the predicted price (e.g., \$450,000).

c. Customer Churn Prediction (Classification)

Businesses want to predict which customers are likely to stop using their service (churn). A model is trained on historical customer data, where each customer is labeled as 'churned' or 'stayed'.

How it's used: Telecom companies, streaming services (like Netflix), and banks use these predictions to proactively offer incentives to at-risk customers to retain them.

Input Features: Customer's tenure, monthly bill, usage data, customer service call history, contract type.

Output: A binary label: churn (1) or no churn (0).

d. Credit Card Fraud Detection (Classification)

This is a critical security application where a model learns to identify fraudulent credit card transactions. The dataset consists of transactions labeled as 'fraudulent' or 'legitimate'.

How it's used: Financial institutions use this in real-time to block potentially fraudulent transactions and alert the cardholder.

Input Features: Transaction amount, time of day, location, merchant category, historical spending patterns.

Output: A binary label: fraud (1) or legitimate (0).

e. Medical Diagnosis (Classification)

Models can be trained to predict the presence of a specific disease based on a patient's medical data.

How it's used: As a tool to assist doctors in diagnosing conditions like diabetes, cancer, or heart disease by identifying patterns in patient data that might not be obvious to the human eye.

Input Features: Patient age, blood pressure, cholesterol levels, results from medical tests (e.g., ECG, blood sugar).

Output: A binary label indicating the presence (1) or absence (0) of the disease.

For this report, we will focus on the Medical Diagnosis application, specifically predicting heart disease.

2. Justification of Results for Heart Disease Prediction

For this task, four different supervised learning models were implemented and evaluated on the UCI Heart Disease dataset. The goal was to classify whether a patient has heart disease based on various medical attributes. The performance was primarily measured by Accuracy, which is the percentage of correct predictions.

Model Performance Summary:

Model	Accuracy on Test Set
Logistic Regression	88.3%
Random Forest	86.7%
Support Vector Machine (SVM)	85.0%
Logistic Regression	83.3%

Analysis of Model Performance

The results clearly indicate that the **K-Nearest Neighbors (KNN)** model performed the best, achieving an accuracy of **88.3%**. The Random Forest model followed closely, while Logistic Regression had the lowest performance among the four.

Why did K-Nearest Neighbors perform the best?

The superior performance of the KNN model on this specific dataset can be attributed to the following key factors:

1. Capturing Local Patterns: KNN is an instance-based learning algorithm. Unlike other models that try to build a general internal model of the data (like a line or a set of tree rules), KNN makes predictions based on the 'k' most similar data points (neighbors) from the training set. This approach is highly effective when the data has distinct local patterns or clusters. In a medical context, this suggests there might be specific profiles of patients where their combined attributes form a close cluster in the feature space, strongly indicating the presence or absence of heart disease. KNN excels at identifying these "neighborhoods."
2. Non-Linear Decision Boundary: The relationship between patient attributes and heart disease is complex and rarely linear. KNN is a non-parametric model, meaning it doesn't assume a specific functional form (like a straight line) for the decision boundary. The boundary it creates is flexible and can adapt to the irregular, complex shapes that often define real-world data, which gives it an advantage over the strictly linear Logistic Regression.
3. Effectiveness of Feature Scaling: The performance of KNN is highly dependent on the distance between data points. In the provided code, StandardScaler was used to normalize the features. This step is crucial because it ensures that features with large numerical ranges (like cholesterol 'chol') do not dominate the distance calculation over features with small ranges (like sex). By scaling the data effectively, we allowed the KNN algorithm to consider each feature's contribution fairly, which was key to its success. While Random Forest is less sensitive to feature scaling, for KNN it is a critical step that unlocked its high performance.

In conclusion, the K-Nearest Neighbors model's ability to capture localized data patterns with a flexible, non-linear boundary, combined with proper feature scaling, made it the most suitable and best-performing model for this heart disease prediction task.