

# Applications of Supervised Learning and Justification of Results

## 1. Introduction

Supervised learning uses labeled data (features → target) to learn predictive mappings for classification or regression tasks. Below are five practical application areas, each with a brief description and where these models are used.

### 1.1) Spam Email Detection (Classification)

Spam email detection is one of the most common applications of supervised learning. A model is trained on labeled emails, some tagged as “spam” and others as “not spam.” By learning patterns in the email content, subject line, or sender information, the model can classify new incoming emails effectively.

**Where these models can be used:** Email providers such as Gmail and Outlook, enterprise email systems, and cybersecurity applications.

### 1.2) House Price Prediction (Regression)

In the real estate industry, supervised learning is applied to estimate property values. Models are trained on historical data containing details such as location, size, number of rooms, and age of the house, with the target variable being the selling price. This helps in predicting the price of new properties.

**Where these models can be used:** Real estate companies, property valuation platforms, banks for mortgage evaluation, and government agencies for taxation purposes.

### 1.3) Disease Prediction (Classification)

Supervised learning is used in healthcare to predict the likelihood of a disease by analyzing patient information such as age, blood pressure, blood sugar levels, and medical history. By training on labeled medical datasets, these models can assist in early diagnosis and risk detection.

**Where these models can be used:** Hospitals, diagnostic centers, health insurance companies, and public health monitoring systems.

### 1.4) Customer Churn Prediction (Classification)

Businesses use supervised learning to identify customers who are likely to discontinue their services. By studying customer behavior, transaction history, and engagement levels, the models can classify whether a customer is likely to churn or remain loyal.

**Where these models can be used:** Telecom industries, online subscription services, banks, e-commerce platforms, and SaaS businesses.

## 1.5) Credit Card Fraud Detection (Classification)

Financial institutions use supervised learning to detect fraudulent transactions. By analyzing transaction history labeled as either legitimate or fraudulent, the models learn to identify unusual or suspicious activity that may indicate fraud.

**Where these models can be used:** Banks, payment gateways, credit card companies, and e-commerce transaction systems.

## Selected Application: House Price Prediction

Among the above applications, the chosen problem for further implementation is **House Price Prediction**. It is a well-defined regression task with widely available datasets and interpretable features, making it suitable for comparing multiple supervised learning models.

- Clean, well-understood regression problem suitable for comparing diverse models.
- Rich, interpretable features (location, area, rooms, age) make it easier to justify results.

# 2. Justification of Results

In this project, four supervised learning models—Linear Regression, Decision Tree, Random Forest, and Gradient Boosting—were applied to the Ames Housing dataset to predict house prices. Their performance was compared using metrics such as MAE, RMSE, and  $R^2$  score, which measure prediction accuracy and explanatory power.

## 2.1 Dataset Characteristics

The dataset contains both numerical and categorical variables including lot area, number of rooms, year built, and neighborhood. The relationship between these features and housing prices is often non-linear and influenced by feature interactions. Therefore, models need to handle complexity, outliers, and mixed data types to perform well.

## 2.2 Model-wise Analysis

### Linear Regression

- Assumes a straight-line relationship between features and target.
- It underperformed since house prices depend on non-linear factors (e.g., quality or neighborhood influence).

- While simple and interpretable, it was unable to capture interactions, leading to **underfitting**.

### **Decision Tree**

- Captured non-linear relationships better than Linear Regression.
- However, a single tree tended to **overfit**, making predictions unstable.
- Performance was decent, but it lacked robustness for unseen data.

### **Random Forest**

- An ensemble of multiple decision trees, averaging their outputs.
- Reduced overfitting and handled both categorical and numerical variables effectively.
- Delivered much higher accuracy and stability, showing why ensemble learning works well for structured datasets.

### **Gradient Boosting**

- Builds trees sequentially, where each new tree corrects the previous one's mistakes.
- Provided the highest accuracy among all models, as it could capture subtle patterns in the data.
- However, it required careful parameter tuning and was more computationally expensive.

## **2.3 Why Ensembles Performed Best**

Both Random Forest and Gradient Boosting outperformed the simpler models because they:

- Handled complex non-linear interactions effectively.
- Reduced variance and overfitting compared to a single Decision Tree.
- Balanced bias and variance, leading to stronger generalization.

## **2.4 Final Justification**

The results are justified by the dataset's complexity. Linear Regression was too simplistic, and Decision Tree was prone to instability. Random Forest improved generalization by averaging multiple trees, while Gradient Boosting achieved the best accuracy by sequentially minimizing errors.

Hence, ensemble models proved most suitable for house price prediction, as they combined predictive strength, robustness, and adaptability to real-world data patterns.