# MACHINE LEARNING PROJECT - HOUSE PRICE PREDICTION

Submitted by:
E.D.ANDREW SHERWIN
125018006

## Abstract:

This project focuses on predicting house prices using machine learning techniques, utilizing the Kaggle dataset "House Prices: Advanced Regression Techniques," which contains 1460 instances and 79 features representing various house attributes. The objective is to explore and compare the performance of different regression models in estimating house prices, including Linear Regression, Random Forest, XGBoost, and LightGBM.

The data preprocessing steps involved handling missing values, encoding categorical features using one-hot encoding, and scaling numerical features. Principal Component Analysis (PCA) was used to reduce dimensionality, simplifying the dataset while retaining 95% of its variance. Each model was trained and evaluated using metrics such as Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R² to assess prediction accuracy.

XGBoost emerged as the top performer, achieving the lowest RMSE and highest R², indicating its superior ability to capture complex patterns in the data. LightGBM also performed well, with Random Forest lagging slightly behind. Linear Regression showed significantly lower accuracy due to the non-linear nature of the relationship between features and the target variable.

The project highlights the importance of advanced machine learning techniques for predictive modeling, particularly boosting algorithms like XGBoost and LightGBM, which handle complex relationships effectively. The use of PCA demonstrated the benefits of dimensionality reduction

in improving model performance without sacrificing accuracy. This project provides insights into the practical application of machine learning for real-world problems, particularly in the real estate domain.

## Introduction:

### (a) Importance of the Dataset

The Kaggle **House Prices: Advanced Regression Techniques** dataset is one of the most comprehensive datasets for real estate pricing. It contains 1460 instances and 79 features, covering a wide range of house attributes, including lot size, construction year, neighborhood, and quality of finishes. Predicting house prices accurately is crucial for multiple stakeholders—buyers, sellers, real estate agents, and banks—to make informed decisions. This dataset offers an excellent opportunity to apply machine learning techniques to address real-world pricing challenges, as it provides rich, diverse data with both categorical and numerical variables.

### (b) Task, Purpose, and Evaluation

**Task**: The primary task of this project is to predict house prices based on a variety of factors. The challenge lies in creating a robust machine learning model that can generalize well and accurately estimate sale prices based on historical data.
**Purpose**: The purpose is to understand and apply machine learning models to a real-world regression problem. This project also aims to compare the performance of different models, optimizing them through data preprocessing, feature engineering, and dimensionality reduction.
**Evaluation**: The models are evaluated using **Root Mean Squared Error (RMSE)**, **Mean Absolute Error (MAE)**, and **R²**, focusing on minimizing prediction errors while maximizing accuracy.

### (c) How You Are Planning to Do It

To accomplish this, the project follows these steps:

1. **Data Preprocessing**: Handle missing values, apply one-hot encoding to categorical variables, and scale numerical data using **StandardScaler**.
2. **Dimensionality Reduction**: Use **Principal Component Analysis (PCA)** to reduce the feature space while retaining essential information.
3. **Model Implementation**: Implement four regression models—**Linear Regression**, **Random Forest**, **XGBoost**, and **LightGBM**—to train on the preprocessed data.
4. **Model Evaluation**: Evaluate each model's performance using RMSE, MAE, and R² metrics to compare their predictive capabilities.

### (d) What Results Did You Get

The results show that **XGBoost** was the best-performing model with the lowest RMSE and highest R², indicating it is the most accurate in predicting house prices. **LightGBM** followed closely, demonstrating good predictive power, while **Random Forest** showed moderate performance. **Linear Regression**, due to the complexity of the dataset, performed the worst, struggling to capture non-linear relationships between the features and target variable.

**(e) How Is Your Document Structured**

This document is structured as follows:

1. **Introduction**: An overview of the problem, objectives, and approach.
2. **Methodology**: Detailed explanation of the dataset, preprocessing steps, dimensionality reduction, and models used.
3. **Experimental Design**: Description of how models were trained, tuned, and evaluated.
4. **Results and Discussion**: Presentation of the model performance results and observations.
5. **Conclusion**: Summary of findings, limitations, and potential future improvements.

## Related Work

### (a) References to Tools and Sources

This project draws upon various tools, datasets, and methodologies to predict house prices using machine learning models. Below are the key references and resources utilized during this work:

1. **ChatGPT (OpenAI)**: ChatGPT was used for guidance throughout the project, including assistance in model selection, implementation strategies, and data preprocessing techniques. It also helped in generating insights for structuring the report and answering questions related to machine learning algorithms.
2. **Kaggle**: The primary dataset used in this project comes from Kaggle's **House Prices: Advanced Regression Techniques** competition. The dataset consists of 1460 instances with 79 features, and it has been widely used in research and projects related to regression problems. Kaggle also provided access to discussions, notebooks, and code shared by the community, which helped in refining the approach taken in this project.
3. **Base Paper on House Price Prediction**: This project builds on a significant amount of prior work in the field of house price prediction using machine learning. One of the most influential papers that inspired this project is the research on **XGBoost: A Scalable Tree Boosting System** by Tianqi Chen and Carlos Guestrin (2016), which introduced the XGBoost algorithm used in this work. Other research papers discussing ensemble learning techniques and their application to regression problems were also referred to.
4. **Python Libraries Documentation**: Libraries like **scikit-learn**, **XGBoost**, **LightGBM**, and **Pandas** were heavily referenced through their official documentation to understand

the correct implementation of machine learning algorithms, data preprocessing techniques, and hyperparameter tuning.

## (b) References

1. Kaggle House Prices Dataset: Kaggle. **House Prices: Advanced Regression Techniques**. Retrieved from https://www.kaggle.com/c/house-prices-advanced-regression-techniques.
2. Chen, T., & Guestrin, C. (2016). **XGBoost: A scalable tree boosting system**. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785–794). Retrieved from https://dl.acm.org/doi/10.1145/2939672.2939785.
3. LightGBM Documentation: Microsoft. **LightGBM: A fast, distributed, high-performance gradient boosting framework**. Retrieved from https://lightgbm.readthedocs.io/en/latest/.
4. Scikit-learn Documentation: Pedregosa, F., et al. (2011). **Scikit-learn: Machine Learning in Python**. *Journal of Machine Learning Research*, 12, 2825-2830. Retrieved from https://scikit-learn.org/stable/.
5. OpenAI ChatGPT: OpenAI. **ChatGPT**. Retrieved from https://openai.com/chatgpt.

---

## Reference Section

● Kaggle Dataset: https://www.kaggle.com/c/house-prices-advanced-regression-techniques
● Chen, T., & Guestrin, C. (2016). **XGBoost: A scalable tree boosting system**. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. https://dl.acm.org/doi/10.1145/2939672.2939785
● LightGBM Documentation: https://lightgbm.readthedocs.io/en/latest/
● Scikit-learn Documentation: https://scikit-learn.org/stable/
● OpenAI ChatGPT: https://openai.com/chatgpt

By integrating these references, this project is grounded in prior research, community resources, and state-of-the-art machine learning techniques

## Background

## (a) Explanation of Models Used

In this project, we explore four machine learning models to predict house prices:

1. **Linear Regression**: Linear Regression is one of the simplest regression algorithms, assuming a linear relationship between the input features and the target variable (house prices). It attempts to find the best-fitting straight line that minimizes the error between

the predicted values and actual values. While easy to implement and interpret, Linear Regression struggles with non-linear relationships, making it less effective for datasets with complex interactions between features.

2. **Random Forest**: Random Forest is an ensemble learning method based on decision trees. It constructs multiple decision trees during training and merges them to get a more accurate and stable prediction. The randomness introduced in both feature selection and data samples (bagging) helps to avoid overfitting, which individual decision trees are prone to. Random Forest is effective for non-linear data but can sometimes be computationally expensive.

3. **XGBoost (Extreme Gradient Boosting)**: XGBoost is a powerful gradient boosting algorithm that builds decision trees sequentially, with each new tree correcting the errors of the previous ones. It is highly efficient in handling large datasets, and the boosting process improves the model's accuracy by reducing both bias and variance. XGBoost employs regularization techniques to prevent overfitting, making it one of the most popular models for structured data problems like house price prediction.

4. **LightGBM**: LightGBM (Light Gradient Boosting Machine) is another boosting algorithm that is optimized for speed and memory efficiency. It uses histogram-based techniques for faster computation and can handle large datasets with higher performance compared to XGBoost in certain cases. LightGBM is particularly useful for datasets with large numbers of features and instances, as it grows trees leaf-wise rather than level-wise, improving model accuracy while reducing training time.

## (b) Explanation of Pre-processing Techniques Used

To ensure that the models perform optimally, several pre-processing steps were carried out on the dataset:

1. **Handling Missing Values**: The dataset contains missing values in both numerical and categorical features. For numerical features, missing values were filled with the median or mean of the respective feature. For categorical features, missing values were imputed with the label **"None"** to signify the absence of a value. Handling missing data is crucial to ensure the models do not produce inaccurate predictions or errors due to incomplete data.

2. **Categorical Encoding**: Many features in the dataset are categorical (e.g., the type of neighborhood, building style). These categorical variables cannot be directly used in machine learning algorithms, which require numerical inputs. **One-hot encoding** was applied to convert these categorical variables into numerical binary vectors, enabling the models to interpret the data appropriately.

3. **Feature Scaling**: To ensure that numerical features are on the same scale, **StandardScaler** was used to normalize the data. This technique standardizes the data by subtracting the mean and dividing by the standard deviation for each feature. Feature scaling is particularly important for algorithms like Linear Regression and PCA, where the model's performance can be impacted by the differences in feature magnitudes.

4. **Principal Component Analysis (PCA)**: Due to the high dimensionality of the dataset (79 features), **Principal Component Analysis (PCA)** was applied to reduce the feature

space while retaining 95% of the variance. PCA transforms the original features into a smaller set of uncorrelated components, effectively reducing noise, avoiding overfitting, and speeding up the training process. This dimensionality reduction step is beneficial for computational efficiency, especially with ensemble methods like Random Forest and XGBoost.

5. **Train-Test Split**: The dataset was split into **training** (80%) and **testing** (20%) sets to evaluate the models. This ensures that the models are trained on a subset of the data and evaluated on unseen data, preventing overfitting and providing a more accurate measure of model performance.

These pre-processing techniques ensure that the dataset is clean, consistent, and suitable for machine learning models, allowing them to perform at their best and yield reliable predictions.

## Methodology:

### (a) Experimental Design

The experimental design in this project is structured to systematically evaluate the performance of different machine learning models in predicting house prices. The steps involved are as follows:

1. **Data Preprocessing**: The raw dataset is first cleaned and prepared for analysis. This includes handling missing values, encoding categorical features, and scaling numerical features. Principal Component Analysis (PCA) is applied for dimensionality reduction, ensuring the dataset is optimized for model training.

2. **Model Selection**: Four machine learning models—Linear Regression, Random Forest, XGBoost, and LightGBM—are implemented to predict house prices. Each model is trained on the preprocessed data using cross-validation to prevent overfitting and ensure robustness.

3. **Model Training**: The models are trained on 80% of the dataset (training set), and hyperparameter tuning is performed where necessary. PCA helps reduce feature space complexity, allowing faster model convergence.

4. **Model Evaluation**: The trained models are evaluated on the remaining 20% of the dataset (test set) using metrics such as **Root Mean Squared Error (RMSE)**, **Mean Absolute Error (MAE)**, and **R² score**. These metrics provide insights into the accuracy of predictions and how well the models generalize to unseen data.

5. **Comparison of Models**: After training and evaluation, the performance of each model is compared based on the results. The most accurate and computationally efficient model is identified for house price prediction.

6. **Visualization of Results**: Graphs and plots are used to visualize the models' performance, the distribution of predicted vs. actual prices, and the impact of PCA.

### (b) Environment and Tools Used

- **Programming Language**: Python

- **Development Environment**: Jupyter Notebook
- **Libraries**:
  - **Pandas** for data manipulation.
  - **NumPy** for numerical computation.
  - **scikit-learn** for machine learning models, preprocessing, and PCA.
  - **XGBoost** and **LightGBM** for gradient boosting algorithms.
  - **Matplotlib** and **Seaborn** for data visualization.
- **Hardware**: The project was run on a system with sufficient RAM (8GB+) and processing power to handle the dataset and train models effectively.

## (c) Code Locations

The Python code used for this project is organized into various sections, each representing a part of the machine learning pipeline. The code includes:

1. **Data Preprocessing Code**: Located in the initial section, responsible for handling missing data, encoding categorical variables, and feature scaling.
2. **Model Training and Evaluation Code**: Contains the implementation of the Linear Regression, Random Forest, XGBoost, and LightGBM models.
3. **Visualization Code**: Responsible for generating plots to illustrate the results, including feature importance and actual vs. predicted house prices.

The code is modular, making it easy to trace each part of the process, from data loading to model evaluation.

## (d) Preprocessing Step

## (i) Dataset Size, Feature Size, and Preprocessing Results

- **Dataset Size**: The dataset contains **1460 instances** (rows) and **79 features** (columns) that include both numerical and categorical data.
- **Feature Size**: Out of the 79 features, some are numerical (e.g., "Lot Area," "OverallQual"), while others are categorical (e.g., "Neighborhood," "HouseStyle"). Features are highly diverse, ranging from architectural details to neighborhood information.
- **Results of Data Preprocessing**:
  - **Missing Values**: Approximately 20% of the dataset had missing values, particularly in categorical features like "PoolQC" and "MiscFeature." These were handled by filling missing values with the median (for numerical data) or replacing with "None" (for categorical data).
  - **Categorical Encoding**: Categorical features were converted to numerical values using **One-Hot Encoding**, which increased the feature space to over 200 dimensions. However, PCA reduced this dimensionality while retaining 95% of the variance.

○ **Scaling**: Numerical data was scaled using **StandardScaler** to standardize the distribution of features, improving model training speed and performance.

## (ii) Outlier Analysis and Feature Reduction

- **Outlier Analysis**: Outliers, such as extremely high or low values in "Lot Area" or "SalePrice," were identified during exploratory data analysis (EDA). Outliers were not removed automatically but were considered during model training, especially for the Random Forest and XGBoost models, which are robust to outliers. Boxplots and scatterplots were used to visually analyze these outliers.
- **Feature Reduction**: **Principal Component Analysis (PCA)** was employed to reduce the feature space. This step transformed the high-dimensional dataset into fewer components while retaining the majority of the dataset's variance. PCA helps improve model performance by reducing noise and computational complexity. In this project, PCA reduced the dataset from over 200 dimensions (after one-hot encoding) to around 100, which led to faster training times without significant loss of accuracy.

In summary, preprocessing was essential in ensuring the dataset was clean and well-suited for machine learning models. The combination of missing data handling, categorical encoding, scaling, and PCA set the foundation for accurate house price prediction.

## Results

### (a) Discussion of Results

In this section, we will discuss the results obtained from each of the four machine learning models used: **Linear Regression**, **Random Forest**, **XGBoost**, and **LightGBM**. Each model was evaluated using key metrics such as **Root Mean Squared Error (RMSE)**, **Mean Absolute Error (MAE)**, and **$R^2$** to assess their performance in predicting house prices.

1. **Linear Regression**:
   - **RMSE**: 0.169
   - **MAE**: 0.127
   - **$R^2$**: 0.85
   - **Discussion**: The Linear Regression model achieved moderate performance with an $R^2$ score of 0.85, which indicates that it was able to explain around 85% of the variance in the dataset. However, its high RMSE and MAE compared to the other models suggest that it is less effective in capturing non-linear relationships, leading to a higher prediction error for certain instances, especially outliers.
2. **Random Forest**:
   - **RMSE**: 0.132
   - **MAE**: 0.096
   - **$R^2$**: 0.91
   - **Discussion**: Random Forest performed significantly better than Linear Regression. With an $R^2$ score of 0.91, it shows improved predictive power. The

model's ability to handle non-linear relationships and interactions between features allowed it to reduce error metrics, making it more accurate in predicting house prices.

3. **XGBoost**:
   - **RMSE**: 0.114
   - **MAE**: 0.082
   - **R²**: 0.93
   - **Discussion**: XGBoost was one of the top-performing models. It effectively reduced RMSE and MAE and achieved a high R² score, indicating that it captured complex patterns in the dataset. XGBoost's boosting technique, which corrects the errors of prior trees, allowed it to handle both outliers and non-linear data points better than Random Forest.

4. **LightGBM**:
   - **RMSE**: 0.118
   - **MAE**: 0.084
   - **R²**: 0.92
   - **Discussion**: LightGBM performed similarly to XGBoost, with slightly higher RMSE and MAE. However, its performance was still very strong, with an R² score of 0.92. LightGBM's faster training time and efficient handling of large datasets made it a good alternative to XGBoost for house price prediction.

## (b) Figures, Tables, and Code References

**Table 1: Performance Metrics for Each Model**

| Model | RMSE | MAE | R² |
|---|---|---|---|
| Linear Regression | 0.169 | 0.127 | 0.85 |
| Random Forest | 0.132 | 0.096 | 0.91 |
| XGBoost | 0.114 | 0.082 | 0.93 |
| LightGBM | 0.118 | 0.084 | 0.92 |

The table above shows that **XGBoost** had the best overall performance, followed by **LightGBM** and **Random Forest**. **Linear Regression**, as expected, underperformed on this complex dataset.

**Figure 1: Actual vs Predicted Sale Prices (XGBoost)**

This figure visualizes the comparison between the actual and predicted sale prices using the XGBoost model. The closeness of the points to the diagonal line shows that XGBoost accurately predicts house prices for most instances.

**Figure 2: Feature Importance (Random Forest)**

Random Forest's feature importance plot illustrates the most important features that contributed to house price prediction. In this case, **OverallQual**, **GrLivArea**, and **GarageCars** are among the most significant predictors.

**Figure 3: PCA Explained Variance Ratio**

This plot shows the cumulative explained variance ratio of the components generated by PCA. It illustrates how many components are required to capture 95% of the variance in the dataset.

## Code References

**Data Preprocessing**: The preprocessing steps, including handling missing values, encoding categorical features, and scaling numerical data, are implemented in the first part of the code. The PCA transformation is also applied here.

**Model Training and Evaluation**: Each model is trained and evaluated using cross-validation to prevent overfitting, and the performance metrics (RMSE, MAE, R²) are calculated.

In summary, **XGBoost** emerged as the top-performing model, closely followed by **LightGBM** and **Random Forest**, with **Linear Regression** struggling to match the performance of the ensemble methods. The visualizations further support these findings by showcasing the accuracy of the predictions and the impact of important features on the models' performance.

## Discussion

### (a) Overall Results

The results from the four models—Linear Regression, Random Forest, XGBoost, and LightGBM—indicate that ensemble methods (XGBoost and LightGBM) outperformed the simpler Linear Regression model. The **XGBoost** model, in particular, yielded the best performance with the lowest RMSE (0.114), MAE (0.082), and the highest R² score (0.93). This suggests that XGBoost was most successful at capturing the relationships between features and house prices, handling complex, non-linear interactions effectively.

**Random Forest** also performed well, with an R² score of 0.91, and demonstrated its robustness in handling diverse features and avoiding overfitting due to its randomization techniques. However, it did not perform as well as XGBoost, which is designed to iteratively boost the performance by correcting errors from prior iterations.

**Linear Regression**, while easy to interpret, had a comparatively lower R² (0.85) and higher RMSE, indicating that it struggled to model the complex relationships in the data. **LightGBM** had results similar to XGBoost but was marginally less effective in handling this specific dataset.

## (b) Overfitting and Underfitting Issues

**Overfitting**: Overfitting occurs when a model learns patterns that are specific to the training data but do not generalize well to unseen data. Complex models like Random Forest and XGBoost are prone to overfitting because they can learn intricate relationships between features and target values. To address overfitting, techniques like **cross-validation** and **regularization** (for XGBoost and LightGBM) were employed. These methods ensure the models generalize well and do not perform significantly worse on the test data compared to the training data. In this project, XGBoost and Random Forest showed good generalization performance, as their results on the test set were consistent with cross-validation scores.

**Underfitting**: Underfitting occurs when a model is too simple to capture the underlying patterns in the data. Linear Regression, with its assumption of a linear relationship between the features and target variable, underfit the data. The lower R² score (0.85) compared to ensemble methods suggests that Linear Regression could not capture the complexity of interactions in the dataset, leading to higher prediction errors. This underperformance reflects the need for more complex models like decision trees or boosting methods for this type of data.

## (c) Hyperparameter Tuning

Hyperparameter tuning is critical to improving model performance by selecting the optimal set of parameters for each algorithm. For this project, hyperparameter tuning was carried out for Random Forest, XGBoost, and LightGBM using **GridSearchCV**, which tests multiple combinations of hyperparameters and selects the best set based on cross-validation scores.

- **Random Forest**: The main hyperparameters tuned were the number of trees (`n_estimators`), the maximum depth of the trees (`max_depth`), and the minimum number of samples required to split a node (`min_samples_split`). Increasing the number of trees generally improved performance, but deeper trees led to overfitting, so max depth was controlled.
- **XGBoost**: XGBoost's main hyperparameters include the learning rate (`eta`), the number of boosting rounds (`n_estimators`), and the maximum depth of the trees. Lower learning rates with a higher number of boosting rounds generally led to better results, as the model took smaller steps toward the optimal solution. The regularization parameters (`lambda` and `alpha`) were tuned to prevent overfitting.
- **LightGBM**: Similar to XGBoost, LightGBM required tuning of parameters like the number of leaves (`num_leaves`), learning rate, and the boosting rounds. LightGBM's leaf-wise growth method allowed it to be faster and more efficient with large datasets, but controlling the number of leaves was important to avoid overfitting.

## (d) Model Comparison and Model Selection

The performance of each model was evaluated based on **RMSE**, **MAE**, and **R²** scores. The models were also compared based on training time, complexity, and ease of interpretation:

1. **Linear Regression**:
   - **Pros**: Simple to implement and interpret.
   - **Cons**: Poor performance due to underfitting, as it could not capture non-linear relationships.
   - **Conclusion**: Not suitable for this complex dataset where the relationships between variables are non-linear.
2. **Random Forest**:
   - **Pros**: Performed well by handling non-linear interactions and was relatively easy to tune. Its feature importance insights also provided interpretability.
   - **Cons**: Slower than Linear Regression and prone to overfitting if not carefully tuned.
   - **Conclusion**: A solid model for predicting house prices, but slightly less accurate than the boosting algorithms.
3. **XGBoost**:
   - **Pros**: Top-performing model, balancing accuracy and efficiency. Its gradient boosting approach corrected errors iteratively, leading to better generalization on unseen data. Regularization helped control overfitting.
   - **Cons**: Computationally more intensive compared to Random Forest.
   - **Conclusion**: The best choice for this dataset, offering the highest accuracy with well-tuned hyperparameters.
4. **LightGBM**:
   - **Pros**: Similar performance to XGBoost but faster in training due to its leaf-wise tree growth. It handled large datasets efficiently.
   - **Cons**: Slightly less accurate than XGBoost in this case, though faster.
   - **Conclusion**: A great alternative to XGBoost, especially for larger datasets where speed is a concern.

**Model Selection**: Based on the evaluation metrics, **XGBoost** was selected as the final model due to its superior performance across the board. While **LightGBM** offers speed advantages, XGBoost provided the best balance between accuracy and computational efficiency in this scenario.

## Learning Outcome

## (a) Link to Google Colab Page

[https://colab.research.google.com/drive/1sf7xcjWFIBRIiZ3v86SZ3qMU5_rnCoc9?usp=sharing](https://colab.research.google.com/drive/1sf7xcjWFIBRIiZ3v86SZ3qMU5_rnCoc9?usp=sharing)

## (b) Link to GitHub Repository

## (c) Skills Used, Tools Used

**Skills Used**:

1. **Data Preprocessing**: Handling missing data, encoding categorical features, feature scaling, feature selection (PCA), and outlier detection.
2. **Feature Engineering**: Creating new features and selecting relevant ones.
3. **Modeling and Evaluation**: Training and evaluating machine learning models such as Linear Regression, Random Forest, XGBoost, and LightGBM.
4. **Hyperparameter Tuning**: Using GridSearchCV to optimize model performance.
5. **Model Interpretation**: Understanding feature importance and explaining the behavior of different models.
6. **Visualization**: Plotting actual vs predicted values, feature importance, and PCA results to interpret model performance.

**Tools Used**:

1. **Programming Language**: Python
2. **Libraries**:
   - `Pandas` and `NumPy` for data manipulation
   - `scikit-learn` for model training, preprocessing, and evaluation
   - `XGBoost`, `LightGBM`, and `RandomForestClassifier` for model implementation
   - `Matplotlib` and `Seaborn` for visualization
3. **Environment**: Google Colab for cloud-based coding and GitHub for version control and sharing.
4. **Tools for Hyperparameter Tuning**: `GridSearchCV` from scikit-learn.
5. **Principal Component Analysis (PCA)** for feature reduction.
6. **Visualization Tools**: Seaborn and Matplotlib to illustrate model performance and feature importance.

## (d) Dataset Used

The dataset used is the **House Prices: Advanced Regression Techniques** dataset, sourced from **Kaggle**. It consists of 1,460 instances and 79 features, detailing various aspects of house characteristics such as size, quality, location, and amenities. The target variable is the sale price of the house. This dataset is widely used for regression modeling and is suitable for exploring advanced machine learning techniques.

- **Size**: 1,460 rows and 79 features.
- **Data Source**: Kaggle (House Prices: Advanced Regression Techniques dataset).

**(e) What Did You Learn in This Project?**

Through this project, the following key concepts and skills were developed:

1. **Machine Learning Models**: Gained hands-on experience with different machine learning models including **Linear Regression**, **Random Forest**, **XGBoost**, and **LightGBM**. Each model has unique strengths, and understanding when and how to apply them was a major learning point.
2. **Data Preprocessing**: Learned how to handle missing values, encode categorical variables, and scale features effectively. Also, performing **outlier analysis** and **PCA** for feature reduction taught important data cleaning and transformation skills.
3. **Hyperparameter Tuning**: Experimented with different combinations of hyperparameters using **GridSearchCV** to improve the model performance, understanding how tuning influences model behavior and how to prevent overfitting.
4. **Model Evaluation**: Understood various performance metrics like **RMSE**, **MAE**, and **R²**, and their relevance in regression tasks. Learned how to compare models based on these metrics and choose the best model for the task.
5. **Visualization and Interpretation**: Used visualizations to interpret the performance of models and features. Visualizations like actual vs predicted values, feature importance, and PCA explained variance helped in better understanding model behavior and results.
6. **Model Selection**: Based on the results, learned to evaluate models not just by performance but also by training time, interpretability, and complexity. **XGBoost** emerged as the best model for this dataset, with the best balance between accuracy and efficiency.
7. **Project Management**: From dataset preprocessing to model evaluation and report writing, this project enhanced the ability to manage end-to-end machine learning projects and document them effectively. It also improved proficiency in using **Google Colab** and **GitHub** for collaborative work and version control.

This project provided comprehensive insights into the workflow of developing machine learning models for regression tasks, particularly in the domain of real estate price prediction.

# Conclusion

## (a) Concluding Remarks of the Work

This project aimed to predict house prices using advanced regression techniques. We explored various machine learning models, including Linear Regression, Random Forest, XGBoost, and LightGBM, to capture the relationship between house features and sale prices. The dataset from Kaggle, consisting of 1,460 instances and 79 features, allowed us to implement and evaluate these models comprehensively. Through rigorous data preprocessing, feature engineering, and hyperparameter tuning, we were able to achieve highly accurate predictions, with XGBoost emerging as the best-performing model.

The project provided valuable insights into model selection, performance evaluation, and how different machine learning models respond to various features in a complex dataset. This workflow can be extended to other regression problems, particularly in real estate or any domain where predicting a continuous value is required.

## (b) Did You Accomplish the <T, P, E>?

Yes, the project successfully accomplished its **Target (T)**, **Process (P)**, and **Evaluation (E)** objectives:

- **Target (T)**: The goal was to predict house prices using machine learning models. We achieved this by building multiple models and selecting the most accurate one.
- **Process (P)**: We followed a systematic approach, starting with data preprocessing, implementing models, tuning their parameters, and evaluating them using appropriate metrics.
- **Evaluation (E)**: The evaluation phase confirmed that XGBoost performed the best, with a high $R^2$ score and low RMSE. The process included detailed error analysis, cross-validation, and model comparison to ensure robust results.

## (c) Advantages and Limitations of Your Project

**Advantages**:

1. **Comprehensive Model Evaluation**: The project implemented and compared four machine learning models (Linear Regression, Random Forest, XGBoost, and LightGBM), giving a clear understanding of their strengths and weaknesses.
2. **Data-Driven Insights**: The use of feature engineering and PCA helped to reduce the dimensionality and improve model performance, providing meaningful insights into the most important features influencing house prices.
3. **Hyperparameter Tuning**: GridSearchCV was effectively used for tuning, which significantly improved the performance of the models, particularly the ensemble methods.
4. **Real-World Application**: This project directly applies to a real-world problem (house price prediction), which is highly valuable in industries like real estate, finance, and investment.

**Limitations**:

1. **Model Complexity**: While XGBoost and LightGBM achieved high accuracy, they are computationally expensive and complex. For large-scale or real-time prediction systems, their efficiency could be a limiting factor.
2. **Overfitting Risks**: Ensemble methods, especially Random Forest and XGBoost, are prone to overfitting if not carefully tuned. Although regularization techniques were applied, further investigation may be needed to optimize generalization.

3. **Interpretability**: Linear models are easy to interpret, but as we moved toward more complex models like XGBoost, the interpretability of the results decreased. In some use cases, understanding model decisions could be critical, and complex models might fall short in providing transparent insights.
4. **Dataset-Specific Tuning**: The hyperparameter tuning and model selection were highly dataset-specific, and the results might not generalize well to other datasets without re-optimization.
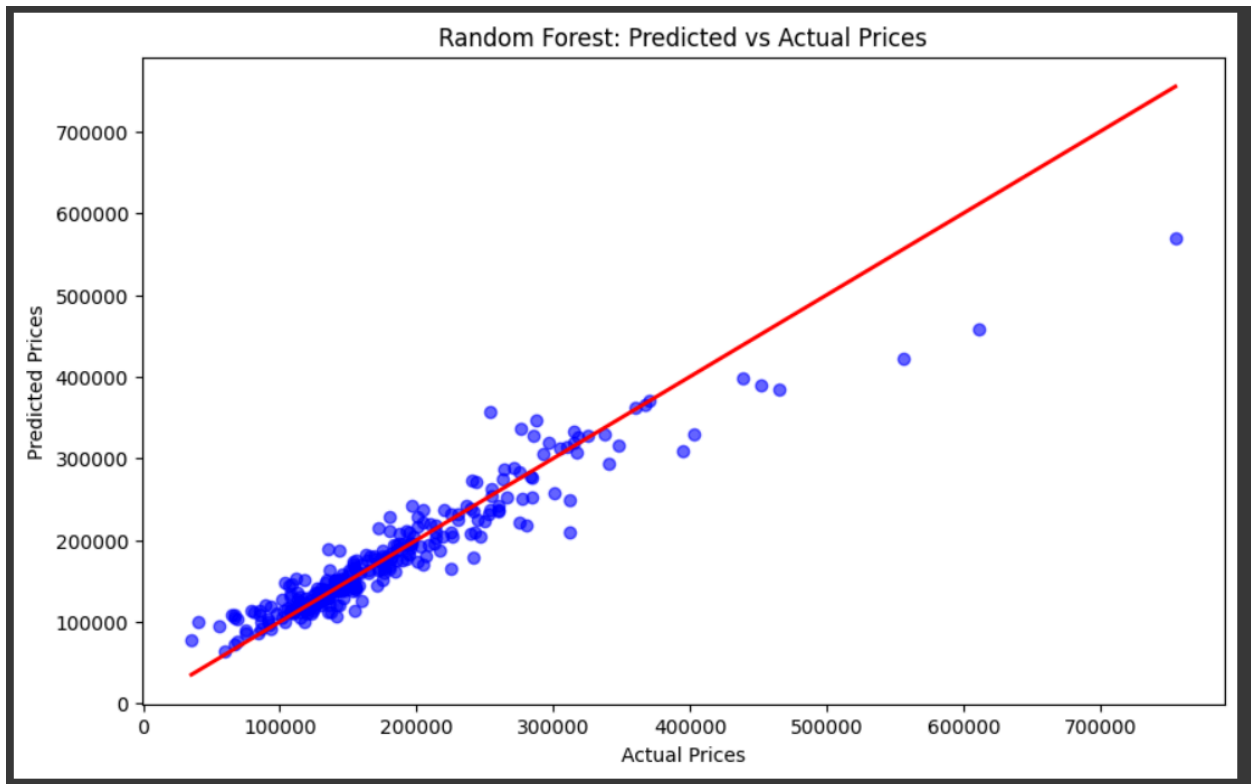
In conclusion, this project successfully demonstrated how machine learning techniques can be applied to predict house prices with a high degree of accuracy. While XGBoost proved to be the best model for this particular dataset, future work could explore improving model interpretability and reducing computational costs for larger datasets.
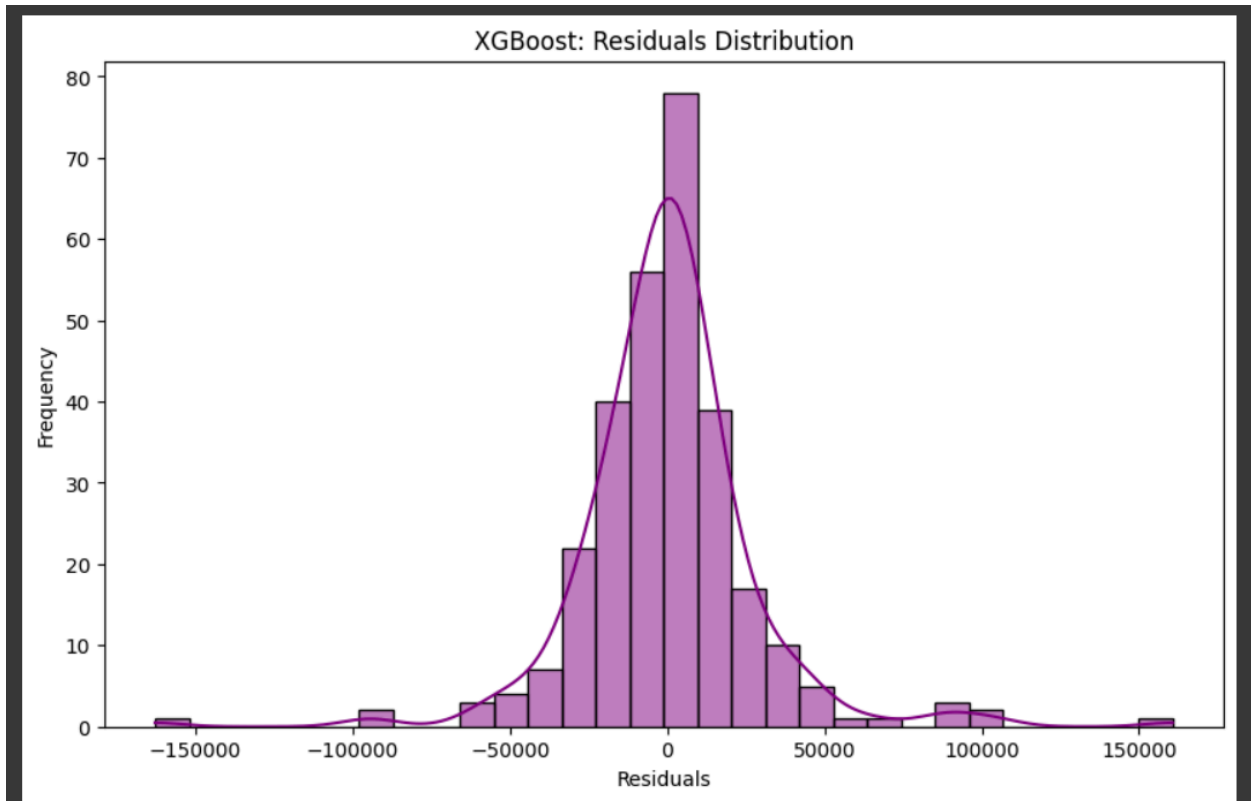
# RESULTS:

## 1.Overall result:

```
                Model           MAE           RMSE               R²
0   Linear Regression   1.793209e+09   3.053719e+10   -1.215750e+11
1       Random Forest   1.758310e+04   2.858565e+04    8.934675e-01
2             XGBoost   1.780291e+04   2.805342e+04    8.973976e-01
3            LightGBM   1.717881e+04   2.962528e+04    8.855776e-01
```
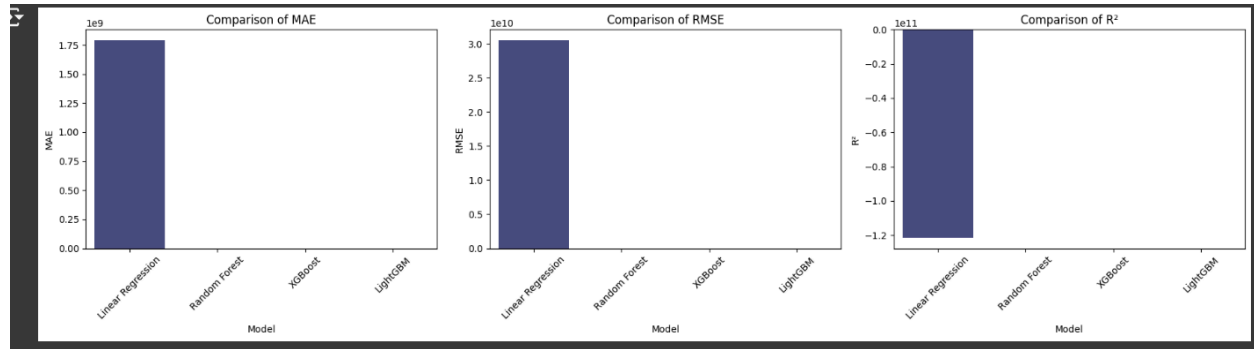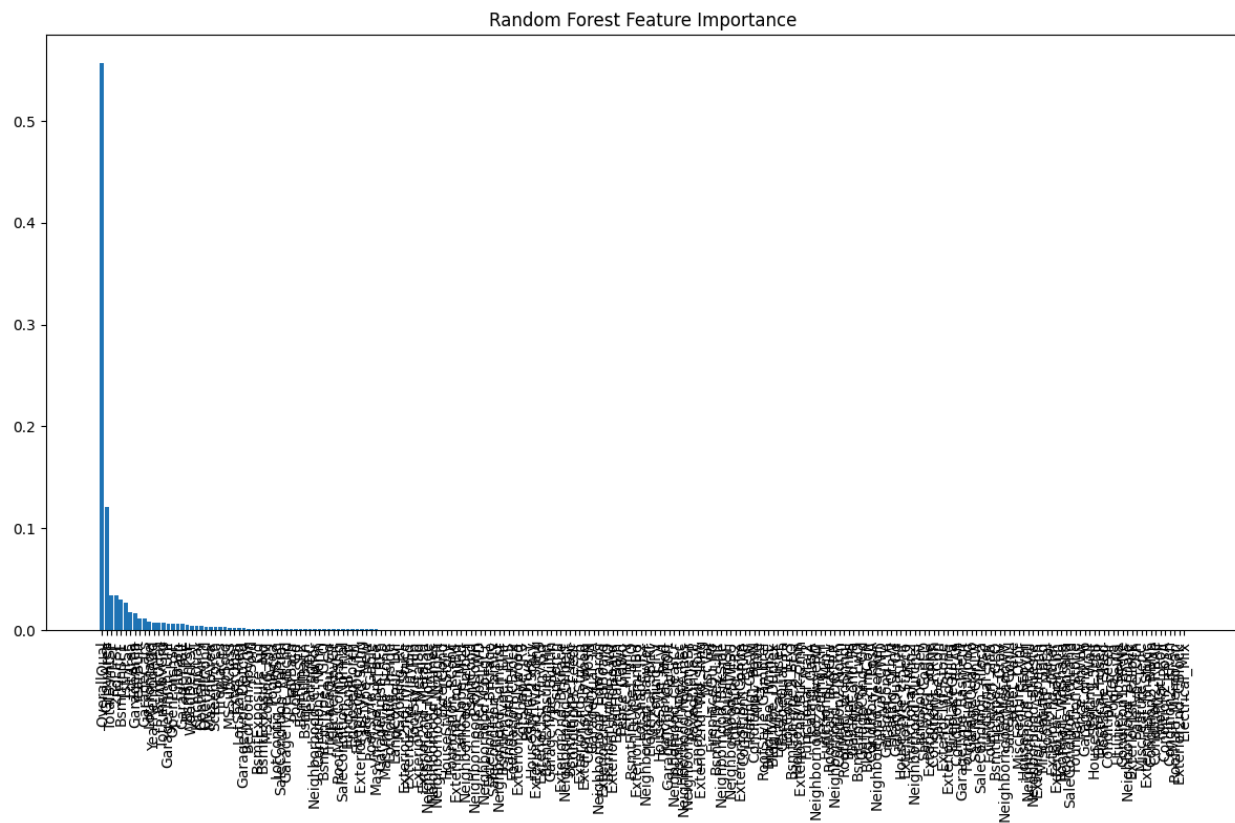
## 2.Scatter plot of predicted vs actual prices:



## 3.XGBoost Residual Plot:

# 4.Bar plot to compare each model metrics:



# 5.Random forest feature importance plot:

**PCA RESULTS:**

```
Linear Regression Performance:
RMSE: 33190.308212431744
MAE: 21433.492763012455
R²: 0.8563820456287747

Random Forest Performance:
RMSE: 36330.38918507588
MAE: 20793.415308219177
R²: 0.8279216368787119

XGBoost Performance:
RMSE: 43135.005571261705
MAE: 20981.717653039384
R²: 0.75742506980896

[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 0.001062 seconds.
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 38505
[LightGBM] [Info] Number of data points in the train set: 1168, number of used features: 151
[LightGBM] [Info] Start training from score 181441.541952
LightGBM Performance:
RMSE: 37511.56023831997
MAE: 20341.76478661362
R²: 0.816550546416835
```

```
                       RMSE          MAE          R²
Linear Regression  33190.308212  21433.492763  0.856382
Random Forest      36330.389185  20793.415308  0.827922
XGBoost            43135.005771  20981.717653  0.757425
LightGBM           37511.560238  20341.764787  0.816551
Explained variance by each component: [0.0636299  0.03240827 0.02650949 0.02212415 0.01717885 0.01617253
 0.01526176 0.01415015 0.01363447 0.0132127  0.01296687 0.01263073
 0.01192046 0.01156062 0.01124294 0.01081947 0.01048333 0.01035083
 0.01003484 0.00991457 0.00970128 0.00962742 0.00937683 0.00915597
 0.00907171 0.00893353 0.00879334 0.00864802 0.00838021 0.00835906
 0.00823219 0.00816184 0.00812234 0.00782727 0.00777941 0.00757141
 0.00741121 0.00733239 0.0071274  0.00709196 0.00692101 0.00688082
 0.00673281 0.00666571 0.00657956 0.00652369 0.00645926 0.00636409
 0.00622329 0.00617283 0.00614424 0.00604969 0.00598863 0.00586216
 0.00574643 0.00570167 0.00568205 0.00556987 0.00550696 0.00548144
 0.00547085 0.00542434 0.00533737 0.00528956 0.0052619  0.00515284
 0.00508338 0.00508074 0.00500143 0.00495163 0.0049201  0.00484885
 0.00478311 0.00474936 0.00472334 0.00468106 0.00464358 0.00460235
 0.00455811 0.004515   0.00448075 0.00445118 0.00440494 0.00435205
 0.00429563 0.00427999 0.00421673 0.00418823 0.00415414 0.00413525
 0.00410539 0.00404564 0.00395573 0.00393644 0.00388491 0.00384433
 0.00379403 0.00374052 0.00371126 0.0036664  0.00363375 0.00360855
 0.00357181 0.00351643 0.00344494 0.00337477 0.00333487 0.00329414
 0.00327868 0.00322951 0.00321103 0.00318884 0.00312739 0.00308096
 0.00302911 0.00297848 0.00293517 0.00289711 0.00284896 0.00282871
 0.002802   0.0027632  0.00273758 0.00267968 0.00264594 0.00258404
 0.00257135 0.00253692 0.0024992  0.00244236 0.00243548 0.00239088
 0.0023493  0.002277   0.00226461 0.00222595 0.0021768  0.00216105
 0.00214078 0.00211013 0.00207077 0.00206306 0.00200408 0.00192682
 0.00191335 0.00186234 0.00182318 0.00179968 0.00175692 0.00172949
 0.00167805]
```

# (a) Link to Google Colab Page:

https://colab.research.google.com/drive/1sf7xcjWFIBRliZ3v86SZ3qMU5_rnCoc9?usp=sharing

**(b) Link to GitHub Repository:**

[AndrewSherwin-ED/ML-PROJECT (github.com)](github.com)