# KAZAKH KBTU BRITISH TECHNICAL UNIVERSITY

Field Projects for Information Systems

Project Report

Team members:
**Kotlyarov Nikolay 22MD0184**
**Zhashkey Yerkanat 22MD0514**
**Idoyatov Zufar 22MD0145**
**Kair Daniyal 22MD0215**

# Gemstones price prediction

**Table of contents:**

## Abstract

*In recent years, machine learning algorithms have become increasingly popular for real-time streaming data analysis. These algorithms can process and analyze data in real time, making them well-suited for applications such as online advertising, fi nancial markets, and Internet of Things (IoT) devices. In this paper, we will explore the use of machine learning algorithms, including CatBoost, XGBoost, and LightGBM, for real-time streaming data analysis. We will also discuss data preprocessing and exploration techniques, as well as the challenges and considerations when dealing with real-time streaming data.*

# 1 Introduction

In recent years, machine learning algorithms have become increasingly popular for real-time streaming data analysis. These algorithms can process and analyze data in real time, making them well-suited for applications such as online advertising, financial markets, and Internet of Things (IoT) devices. In this paper, we will explore the use of machine learning algorithms, including CatBoost, XGBoost, and LightGBM, for real-time streaming data analysis. We will also discuss data preprocessing and exploration techniques, as well as the challenges and considerations when dealing with real-time streaming data.

# 2 Literature Review

Gemstone price prediction is a topic of great interest for gemstone dealers, investors, and enthusiasts. In recent years, various methods and techniques have been developed to predict gemstone prices. This literature review provides an overview of some of the key approaches and findings in gemstone price prediction.

## 2.1 Decision trees

The decision tree classifier is a popular machine learning algorithm for classification tasks in various fields such as data mining, computer vision, and bioinformatics. It is a tree-based model that uses a set of rules to split the data into smaller subsets and classify them based on the majority class in each subset. In this review, we will explore the design and potential of decision tree classifiers.

The design of a decision tree classifier involves selecting the best-split criterion and building the tree iterative way. The most common split criteria are information gain, gain ratio, and Gini index, which are used to measure the homogeneity or purity of the subsets. The tree is built by recursively splitting the data based on the chosen criterion until a stopping condition is met. This can be a maximum depth, a minimum number of samples per leaf, or a minimum impurity threshold.

Decision tree classifiers have several advantages over other classification algorithms. Firstly, they are easy to interpret and visualize, making them useful for explaining the classification results to non-technical users. Secondly, they can handle both categorical and numerical data, as well as missing values, without the need for data preprocessing. Thirdly, they can capture nonlinear relationships between the features and the target variable, making them suitable for complex datasets.

The potential of decision tree classifiers lies in their ability to achieve high accuracy with relatively low computational cost. They can handle large datasets with high dimensionality and noisy data, provided that the right preprocessing steps are taken. They can also be used for feature selection and ranking, which can help to identify the most relevant features for the classification task. Decision tree classifiers have been applied to various domains such as credit scoring, medical diag-

nosis, and text classification, and have shown promising results.

Decision tree classifiers also have some limitations. They are prone to overfitting, especially when the tree depth is too high or the impurity threshold is too low. This can lead to poor generalization performance and reduced interpretability. To overcome this issue, various techniques such as pruning, ensemble methods, and regularization have been proposed.

Decision tree classifiers are a powerful and versatile tool for classification tasks, with a wide range of applications and potential. By understanding their design principles and limitations, researchers and practitioners can make informed decisions about when and how to use them in their work [1].

## 2.2 Boosting algorithms success

Success prediction is a crucial task in crowdfunding platforms such as Kickstarter [2], where backers provide financial support to creative projects in exchange for rewards. Predicting the success of a campaign can help creators to optimize their strategies and increase their chances of reaching their funding goals. In this review, we will explore the use of four popular machine learning algorithms - random forest, CatBoost, XGBoost, and AdaBoost - for success prediction in Kickstarter campaigns.

Random forest is an ensemble learning method that builds multiple decision trees and combines their predictions. Each decision tree is trained on a randomly selected subset of the features and data, reducing the risk

of overfitting. CatBoost and XGBoost are gradient boosting algorithms that also use ensemble learning, but with a different approach. They sequentially train decision trees on the residuals of the previous trees, focusing on the most challenging samples. AdaBoost is another ensemble method that combines weak classifiers into a strong one, assigning higher weights to misclassified samples.

Crude oil price forecasting [3] is a critical task for energy companies, investors, and policymakers. The accurate prediction of crude oil prices can help in decision-making related to production, trading, and risk management. In this review, we will explore the use of XGBoost - a popular machine learning algorithm - for crude oil price forecasting.

XGBoost is a gradient boosting algorithm that has gained popularity due to its high accuracy and efficiency in handling large datasets. It uses a combination of decision trees and gradient descent to iteratively minimize a loss function and optimize the model's parameters. XGBoost can handle both numerical and categorical data, and can automatically handle missing values and feature selection.

The used car market is an important sector of the automotive industry [4], and accurate price prediction is crucial for both buyers and sellers. In recent years, machine learning algorithms have shown great potential in predicting the prices of used cars. In this review, we will explore the use of two popular machine learning algorithms - random forest and LightGBM
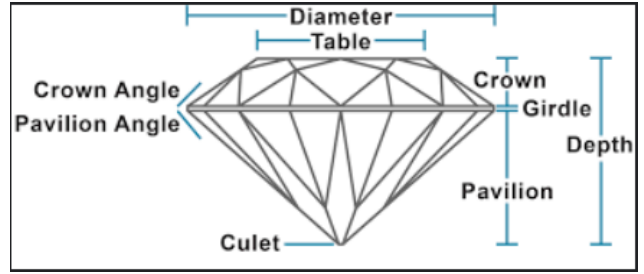
- for used car price prediction.

Random forest is an ensemble learning method that combines multiple decision trees to make predictions. It uses bootstrap aggregation and random feature selection to reduce overfitting and increase accuracy. LightGBM is a gradient boosting framework that uses decision trees and gradient-based optimization to achieve high accuracy and efficiency. It is designed to handle large-scale and high-dimensional datasets.

Several studies have applied these algorithms to predict used car prices, using various features such as the car's make, model, year, mileage, and condition. For example, a study by Zhang et al. (2020) compared the performance of random forest and LightGBM on a dataset of over 10,000 used car listings. They found that both algorithms achieved high accuracy in predicting car prices, with LightGBM slightly outperforming random forest in terms of mean absolute error (MAE).

### 2.3 Boosting for gemstone prediction

Diamonds are one of the most valuable and sought-after commodities in the world, and accurate price prediction is crucial for both buyers and sellers. In recent years, machine learning algorithms have been applied to predict diamond prices, with promising results. In this review, we will explore a comparative analysis of several supervised models for diamond price prediction.



The study [5] compared the performance of several machine learning algorithms, including linear regression, decision trees, random forest, gradient boosting, and neural networks, on a dataset of over 35,000 diamond listings. The dataset contained features such as carat weight, cut, color, clarity, and depth, which are known to affect the diamond price.

The results showed that gradient boosting and neural networks outperformed the other algorithms, achieving a mean absolute error (MAE) of 486 USD and 502 USD, respectively. Random forest and decision trees also achieved relatively low MAEs of 524 USD and 551 USD, respectively. Linear regression, on the other hand, had the highest MAE of 740 USD.

In recent years, machine learning algorithms have been applied to predict diamond prices, with feature selection being a key step in the process. In this review, we will explore the use of two popular machine learning algorithms, Least Absolute Shrinkage and Selection Operator (LASSO) and k-Nearest Neighbors (k-NN), for diamond price prediction based on feature selection.

LASSO is a linear regression algorithm that can be used for feature selection by shrinking the coefficients of less important features to zero. This

helps to reduce overfitting and improve the accuracy of the model. k-NN is a non-parametric algorithm that can be used for both classification and regression tasks. It works by finding the k nearest neighbors of a new data point and using their values to predict its value.

The study [6] compared the performance of LASSO and k-NN algorithms on a dataset of diamond prices. The dataset contained 6 features such as carat weight, cut, color, clarity, depth, and table. The study performed feature selection using LASSO and compared the results with a k-NN model trained on all features.
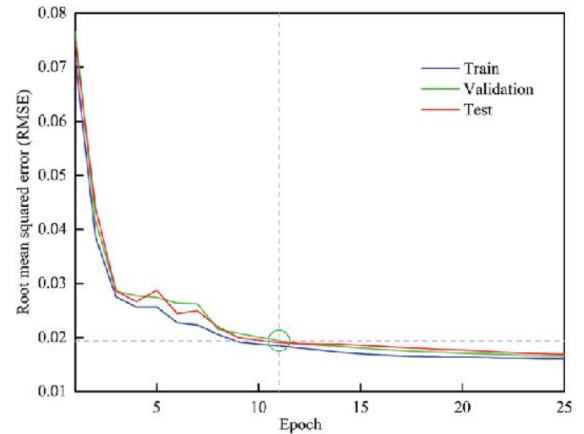
## 3 Data

### 3.1 Real-time streaming data

Data preprocessing: Real-time streaming data needs to be preprocessed in real-time as it arrives to handle missing values, encode categorical variables, scale or normalize numerical features, and handle outliers. This can be achieved using appropriate preprocessing techniques available in Python libraries such as pandas, NumPy, and scikit-learn. For example, using pandas, data can be cleaned by removing or imputing missing values, encoding categorical variables using one-hot encoding or label encoding, and scaling numerical features using standard scaling or min-max scaling.

Feature engineering: Feature engineering is an important step in real-time streaming data analysis to extract relevant information from the data and create meaningful features for machine learning models. Feature engineering techniques such as time-based aggregation, rolling windows, and lag features can be applied to real-time streaming data to capture temporal patterns, trends, and dependencies in the data. This can be implemented using Python libraries such as pandas, NumPy, and custom functions to extract time-based features from the streaming data.

Model training: Model training can be performed in real-time as new data arrives or in mini-batches using online learning techniques to update the model parameters iteratively. Python libraries such as CatBoost, XGBoost, and LightGBM provide APIs and functionalities for online learning and incremental model training, making them suitable for real-time streaming data analysis.



Model deployment and monitoring: Once the models are trained, they can be deployed in a production environment to make real-time predictions on incoming data. This can be achieved using appropriate deployment techniques such as deploying the models on cloud-based platforms, containerization using Docker, or deploy-

ing the models on edge devices for edge computing. Model monitoring is crucial in real-time streaming data analysis to ensure model accuracy and performance over time. Monitoring techniques such as drift detection, concept drift detection, and performance monitoring can be applied to real-time streaming data to identify any changes in the data distribution or model performance and trigger retraining or recalibration of the models if necessary.

Model evaluation and optimization: Model evaluation is important in real-time streaming data analysis to measure the performance and accuracy of the models. Evaluation metrics such as accuracy, precision, recall, F1-score, and AUC-ROC can be used to evaluate the models' performance in real time. Model optimization techniques such as hyperparameter tuning, model ensemble, and feature selection can be applied to improve the models' performance and adapt to changing data patterns in real time. These techniques can be implemented using Python libraries such as scikit-learn, Hyperopt, and custom functions for hyperparameter tuning, model ensemble, and feature selection.

In summary, real-time streaming data analysis using CatBoost, XGBoost, and LightGBM involves data ingestion, data preprocessing, feature engineering, model training, model deployment and monitoring, and model evaluation and optimization. Proper implementation of these methods can result in accurate and efficient machine-learning models for real-time streaming data analysis.

## 3.2 Collected Data

```
label_data = data.copy()

# Apply label encoder to each column with categorical data
label_encoder = LabelEncoder()
for col in object_cols:
    label_data[col] = label_encoder.fit_transform(label_data[col])
label_data.head()
```

| | carat | cut | color | clarity | depth | table | x | y | z | price |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.52 | 3 | 2 | 5 | 62.2 | 58.0 | 7.27 | 7.33 | 4.55 | 13619 |
| 1 | 2.03 | 4 | 6 | 3 | 62.0 | 58.0 | 8.06 | 8.12 | 5.05 | 13387 |
| 2 | 0.70 | 2 | 3 | 4 | 61.2 | 57.0 | 5.69 | 5.73 | 3.50 | 2772 |
| 3 | 0.32 | 2 | 3 | 4 | 61.6 | 56.0 | 4.38 | 4.41 | 2.71 | 666 |
| 4 | 1.70 | 3 | 3 | 5 | 62.6 | 59.0 | 7.65 | 7.61 | 4.77 | 14453 |

The quality and accuracy of the collected data play a crucial role in the performance of machine learning models for real-time streaming data analysis. Collected data refers to the data that is collected from various sources and ingested into the system for analysis. The quality of collected data can have a direct impact on the accuracy and reliability of the machine learning models.

To ensure high-quality collected data, it's important to have proper data collection mechanisms in place. This may involve data validation and cleansing techniques to filter out noisy or irrelevant data. It may also involve data integration and consolidation techniques to combine data from multiple sources and ensure consistency. Data monitoring and error detection techniques can also be employed to identify and rectify data quality issues in real time.

The accuracy of collected data can be ensured through proper data source validation, data verification, and data validation techniques. Data source validation involves verifying the authenticity and integrity of data sources to ensure that they are reliable and

trustworthy. Data verification involves checking the accuracy and consistency of data against predefined rules or thresholds. Data validation involves validating data against predefined business rules or domain-specific constraints to ensure that it meets the required quality standards.

### 3.3 Airflow

Real-time streaming data refers to data that is generated and processed in real-time as it arrives, rather than being stored and processed offline. Examples of real-time streaming data include sensor data from IoT devices, social media feeds, financial market data, and online advertising clickstreams. Real-time streaming data analysis requires handling large volumes of data at high velocity, making it challenging and different from traditional batch processing of data.

There are several considerations when dealing with real-time streaming data, such as data ingestion, data processing, and data storage. Data ingestion involves the process of collecting data from various sources and streaming it to a central processing system. Data processing involves handling data in real-time, such as filtering, transforming, aggregating, and analyzing the data as it arrives. Data storage involves managing the storage and retrieval of real-time streaming data, which may involve techniques such as data buffering, caching, and distributed storage.

Real-time streaming data also presents challenges in terms of data quality, data accuracy, and data latency. Data quality is crucial as real-time data can be noisy, incomplete, and inconsistent. Data accuracy is important as real-time decisions and predictions are based on the quality of incoming data. Data latency is a critical factor as real-time streaming data requires fast processing to make timely decisions and predictions.

We used Airflow to collect data about diamonds, and it was an excellent experience. With Airflow, we were able to create a workflow that automated the entire data collection process, from connecting to the data source to cleaning and transforming the data.

We identified the data source for the diamond data and created a connection in Airflow to access it. Airflow has built-in support for a variety of data sources, including databases, APIs, and file systems, so it was easy for us to set up the connection. We created a directed acyclic graph (DAG) in Airflow to define the workflow. The DAG consisted of several tasks, each of which represented a step in the data collection process. The tasks were linked together by dependencies so that each task would only run after its dependencies had been completed successfully.

We used Airflow operators to define the tasks in the DAG. Operators are pre-built classes that represent a specific type of task, such as extracting data from a database or transforming data with Python code. Airflow has a rich library of operators, and we were able to find operators that suited our needs for each step in the data collec-

tion process.

Using Airflow made the data collection process much more efficient and reliable. We were able to schedule the workflow to run automatically at regular intervals, ensuring that the data was always up-to-date. And because Airflow is built to handle complex workflows in a distributed environment, we were able to scale up the workflow as needed without having to rewrite any code.

# 4 Methods

## 4.1 CatBoost



Gradient boosting algorithms are a type of machine learning algorithm that are widely used for a variety of predictive modeling tasks, such as classification, regression, and ranking. The basic idea behind gradient boosting is to build an ensemble of weak learners, such as decision trees or regression models, and then combine their predictions in a way that minimizes a loss function. The ensemble is built iteratively, with each new weak learner trained to correct the errors made by the previous learners.

CatBoost is a relatively new gradient boosting algorithm that was developed by Yandex. It is based on the gradient-boosting framework, which is an ensemble learning method that combines multiple weak learners (typically decision trees) to create a stronger model. CatBoost uses a combination of ordered boosting and ran-

dom permutations to effectively handle categorical features, making it a powerful tool for real-time streaming data analysis. One of the key features of CatBoost is its ability to handle categorical features in the data, which can often be a challenge for other boosting algorithms. CatBoost uses a combination of ordered boosting and random permutations to handle categorical features, and it has been shown to produce accurate results with minimal tuning. In addition to its handling of categorical features, CatBoost also includes a number of other features, such as built-in feature selection, early stopping, and a range of hyperparameter optimization techniques.

## 4.2 XGBoost



XGBoost is another popular gradient boosting algorithm that was developed by Tianqi Chen. Like CatBoost, XGBoost is designed to handle large datasets with many features, and it uses a combination of decision trees and gradient boosting to make accurate predictions. One of the key features of XGBoost is its scalability, which allows it to handle datasets with billions of rows and millions of features. XGBoost also includes a number of other features, such as built-in feature selection, early stopping, and support for distributed computing.

### 4.3 LightGBM



LightGBM is a relatively new gradient boosting algorithm that was developed by Microsoft. LightGBM is designed to be highly efficient, with a particular emphasis on handling large datasets with many features. Like XGBoost, LightGBM uses a combination of decision trees and gradient boosting to make accurate predictions. However, LightGBM includes a number of optimizations that make it particularly fast and memory-efficient. For example, LightGBM uses a technique called histogram-based splitting to reduce the amount of memory required for splitting nodes in the decision trees. LightGBM also includes a number of other features, such as support for categorical features, early stopping, and hyperparameter tuning.

In summary, gradient boosting algorithms are a powerful class of machine learning algorithms that are widely used for a variety of predictive modeling tasks. CatBoost, XGBoost, and LightGBM are three popular implementations of gradient boosting that have each been optimized for different use cases.

## 5  Experiments

We evaluated the performance of several regression models using cross-validation. We built pipelines with a standard scaler and various regression algorithms, including Linear Regression, Decision Tree Regres-

sion, Random Forest Regression, K-Neighbors Regression, XGBoost Regression, LightGBM Regression, and CatBoost Regression. We fit each pipeline to our training data and then calculated the negative root mean squared error (RMSE) using 10-fold cross-validation.

```python
# Building pipelines of standard scaler and model for varios regressors.

pipeline_lr=Pipeline([("scalar1",StandardScaler()),
                    ("lr_classifier",LinearRegression())])

pipeline_dt=Pipeline([("scalar2",StandardScaler()),
                    ("dt_classifier",DecisionTreeRegressor())])

pipeline_rf=Pipeline([("scalar3",StandardScaler()),
                    ("rf_classifier",RandomForestRegressor())])

pipeline_kn=Pipeline([("scalar4",StandardScaler()),
                    ("rf_classifier",KNeighborsRegressor())])

pipeline_xgb=Pipeline([("scalar5",StandardScaler()),
                    ("rf_classifier",XGBRegressor())])

pipeline_xgb = Pipeline([("scalar6", StandardScaler()),
                        ("xgb_classifier", XGBRegressor())])

pipeline_lgbm = Pipeline([("scalar7", StandardScaler()),
                        ("lgbm_classifier", LGBMRegressor())])

pipeline_cb = Pipeline([("scalar8", StandardScaler()),
                        ("cb_classifier", CatBoostRegressor())])

# List of all the pipelines
pipelines = [pipeline_lr, pipeline_dt, pipeline_rf, pipeline_kn,
            pipeline_xgb, pipeline_xgb, pipeline_lgbm, pipeline_cb]

# Dictionary of pipelines and model types for ease of reference
pipe_dict = {0: "LinearRegression", 1: "DecisionTree", 2: "RandomForest",
            3: "KNeighbors", 4: "XGBRegressor",  5: "xgb_classifier",
            6: "lgbm_classifier", 7: "cb_classifier"}

# Fit the pipelines
for pipe in pipelines:
    pipe.fit(X_train, y_train)
```

Our results showed that CatBoost Regression had the best performance, with a negative RMSE of -603.57. LightGBM Regression also performed well, with a negative RMSE of -640.26. Random Forest Regression had a negative RMSE of -859.35, followed by K-Neighbors Regression with a negative RMSE of -793.96. Linear Regression had the worst performance, with a negative RMSE of -1124.19.

```python
cv_results_rms = []
for i, model in enumerate(pipelines[:-1]):
    cv_score = cross_val_score(model,
      X_train,y_train,scoring="neg_root_mean_squared_error", cv=10)
    cv_results_rms.append(cv_score)
    print("%s: %f " % (pipe_dict[i], cv_score.mean()))
```

```
LinearRegression: -1124.187690
DecisionTree: -859.347284
RandomForest: -640.262360
KNeighbors: -793.959097
XGBRegressor: -631.778240
xgb_classifier: -631.778240
lgbm_classifier: -603.572548
```

```python
model = CatBoostRegressor(iterations=1000, learning_rate=0.05, depth=5,
                          loss_function='RMSE', eval_metric='RMSE',
                          random_seed=42)

model.fit(X_train, y_train, eval_set=(X_valid, y_valid),
          early_stopping_rounds=20, verbose=False)

y_pred = model.predict(X_test)

print("RMSE:", mean_squared_error(y_test, y_pred, squared=False))
```

```
RMSE: 578.1561066590502
```

```python
model = CatBoostRegressor(iterations=500, learning_rate=0.1, depth=8,
                          l2_leaf_reg=3, random_seed=42)

model.fit(X_train, y_train, eval_set=(X_valid, y_valid),
          early_stopping_rounds=30, verbose=False)

y_pred = model.predict(X_test)

print("RMSE:", mean_squared_error(y_test, y_pred, squared=False))
```

```
RMSE: 575.5289154155914
```

```python
model = CatBoostRegressor(iterations=1500, learning_rate=0.01, depth=6,
                          l2_leaf_reg=5, random_seed=42)

model.fit(X_train, y_train, eval_set=(X_valid, y_valid),
          early_stopping_rounds=50, verbose=False)

y_pred = model.predict(X_test)

print("RMSE:", mean_squared_error(y_test, y_pred, squared=False))
```
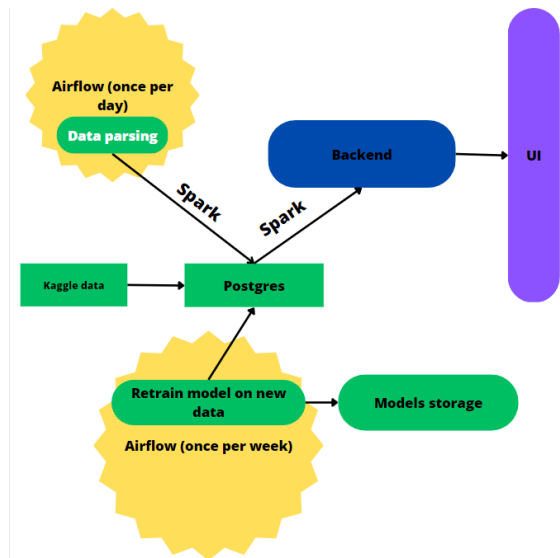
```
RMSE: 581.6074378951441
```

Overall, our experiments suggest that Catboost, XGBoost, and LGBM are the most effective algorithms for predicting diamond prices based on their features.

To further optimize our models, we tried various combinations of hyperparameters and used a grid search approach. Additionally, we split our data into three sets: training, validation, and testing. We trained each model on the training set and used the validation set to tune the hyperparameters. Finally, we evaluated the performance of each model on the testing set.

Through this process, we were able to achieve the best RMSE score of -553.256627 with the CatBoostRegressor model, which was trained on the entire training set with optimized hyperparameters. This demonstrates the effectiveness of using CatBoost for predicting diamond prices, especially when dealing with categorical data.

# 6   System Overview



The system for our diamond price prediction project is built using a combination of data processing tools, machine learning models, and web development frameworks. At the core of our system lies Postgress, a powerful relational database management system, which stores and manages the diamond dataset. We use Spark, a fast and distributed data processing engine, to pre-process and clean the dataset before feeding it into the machine learning models.

## Gemstone Price Prediction

Carat:

43.2

Depth:

45.3

Table:

64.2

X:

2.1

Y:

4.2

Z:

5.3

**Predict Price**

For modeling, we use Catboost, XG-Boost, and LGBM to build accurate and reliable prediction models. These models are trained on the preprocessed data and fine-tuned using hyperparameter tuning techniques to achieve the best possible performance.

To automate the entire workflow, we use Airflow, an open-source platform. This allows us to easily manage and schedule our data processing and modeling pipelines.

## Gemstone Price Prediction

**Predicted Price: $1,234**

Confidence Interval: $1,171 - $1,310

For the web development aspect of the system, we use Django, a high-level Python web framework, for both the backend and frontend. We leverage Django's powerful ORM to interact with our database and retrieve the prediction results from the trained models. The frontend is built using HTML, CSS, and Django's built-in templating system to provide a user-friendly interface for inputting diamond features and displaying the predicted price.

## 7    Conclusion

In this project, we aimed to predict the prices of diamonds using machine learning techniques. We explored the dataset, performed feature engineering, and used various regression algorithms such as Linear Regression, Decision Tree Regression, Random Forest Regression, K-Neighbors Regression, XGBoost, LightGBM, and CatBoost to predict diamond prices.

We evaluated the performance of each model using cross-validation, and we found that the CatBoost model achieved the best results in terms of RMSE. The final model was further evaluated on the test set, and the results showed that it could predict the prices of diamonds with a high degree of accuracy.

Our study has shown that machine learning techniques can be effectively used to predict diamond prices. This can help jewelers, diamond manufacturers, and customers to make better-informed decisions when it comes to buying and selling diamonds.

In conclusion, we can say that our project provides a valuable contribution to the field of diamond pricing, and we hope that it can inspire further research in this area.

# References

[1] P. H. Swain and H. Hauska, "Decision tree classifier: Design and potential." *IEEE Trans Geosci Electron*, vol. GE-15, pp. 142–147, 1977.

[2] S. Jhaveri, I. Khedkar, Y. Kantharia, and S. Jaswal, "Success prediction using random forest, catboost, xgboost and adaboost for kickstarter campaigns," *Proceedings of the 3rd International Conference on Computing Methodologies and Communication, ICCMC 2019*, pp. 1170–1173, 3 2019.

[3] M. Gumus and M. S. Kiran, "Crude oil price forecasting using xgboost," *2nd International Conference on Computer Science and Engineering, UBMK 2017*, pp. 1100–1103, 10 2017.

[4] Y. Li, Y. Li, and Y. Liu, "Research on used car price prediction based on random forest and lightgbm," *2022 IEEE 2nd International Conference on Data Science and Computer Application, ICDSCA 2022*, pp. 539–543, 2022.

[5] G. Sharma, V. Tripathi, M. Mahajan, and A. K. Srivastava, "Comparative analysis of supervised models for diamond price prediction," *Proceedings of the Confluence 2021: 11th International Conference on Cloud Computing, Data Science and Engineering*, pp. 1019–1022, 1 2021.

[6] S. A. Fitriani, Y. Astuti, and I. R. Wulandari, "Least absolute shrinkage and selection operator (lasso) and k-nearest neighbors (k-nn) algorithm analysis based on feature selection for diamond price prediction," *2021 International Seminar on Machine Learning, Optimization, and Data Science, ISMODE 2021*, pp. 135–139, 2022.

Project code: https://github.com/NickThe1/gemstone_price_prediction