Name: ABDULRAHMAN ALMYMAN Student Number

1.IDE:

Jupyter notebook(easy to use) & Google Collab for working online.

2. libraries used to implement the models:

- -Pandas: to read the data easily.
- -Sklearn(include: sklearn.preprocessing, sklearn.model_selection, sklearn.metrics, sklearn.ensemble, etc.)

3. Description of the chosen algorithms including all the selected hyperparameters:

SVR() because it has the best score

Selected Hyperparameters:

default hyperparameters

4. Description of Data Preprocessing Applied:

HW3:

- agardata 2
- 1- Separate features and target variable
- 2- Encode categorical variables with LabelEncoder
- 3- # Initialize the KFold object with shuffle and random_state
- # class sklearn.model_selection.KFold(, *, shuffle=true, random_state=24)
- 4- # Initialize a list to store the mean squared errors, mean absolute errors, RMSE, and R2 scores for each model
- 5- Define the models to evaluate
- 6 Iterate over the models and perform cross-validation & Initialize lists to store the evaluation scores for each fold
- 7 # Perform cross-validation & Split the data into training and testing sets for this fold
- 8 # Fit the model on the training data
- 9-# Make predictions on the testing data
- 10 # Calculate the evaluation scores for this fold
- 11 # Append the evaluation scores to the respective lists

- 12 Calculate the mean of the evaluation scores for this model
- 13 # Append the average evaluation scores to the list
- 14 # Sort the models based on their average mean squared error in ascending order
- 15 # Get the best performing model and its corresponding evaluation scores
- 16- # Print the best performing model and its evaluation scores merging solution(not yet):
- -aqardata_2
- -DocRealestateSale
- -DocRealestateSale2023Q3
- -Transactionssaleforrealestate

5. Description of all other steps that have done to get a final result:

- Selected relevant features for training the model:[mainlocation sublocation neighborhood frontage purpose streetwidth size
 -] and make the Target Pricepm
 - After the all steps we choose svr because it has highest score.

6. Evaluation results:

The results for HW3:

```
Best Model: Support Vector Regression

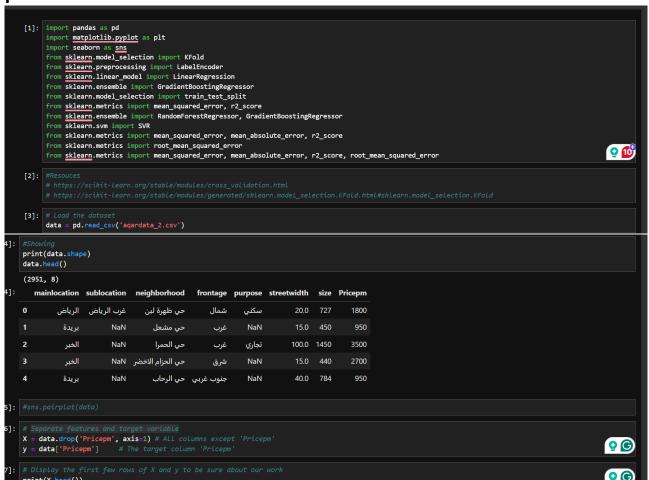
Mean Squared Error: 7081681.56

Mean Absolute Error: 1338.34

Root Mean Squared Error: 2661.14

R<sup>2</sup> Score: 0.00%
```

12. Screenshots of all the code & evaluation results from the platform was chosen:



```
X = data.drop('Pricepm', axis=1) # All columns except 'Pricepm'
y = data['Pricepm'] # The target column 'Pricepm'
                                                                                                                                                                                     Q G
    print(X.head())
                             mainlocation sublocation
              20.0
               بريدة
الخبر
        727
450
        1450
         440
        784
3]: print(y.head())
         1800
           950
         Name: Pricepm, dtype: int64
        encoder = LabelEncoder()
                                                                                                                                                                                       ♀ⓒ
        X_encoded = X.apply(encoder.fit_transform)
[10]: # Initialize the KFold object with shuffle and random_state # class sklearn.model_selection.KFold(n_splits=5, *, shuffle=False, random_state=None)
                                                                                                                                                                                        2
        kf = KFold( shuffle=True, random_state=24)
[11]: # Initialize a list to store the mean squared errors, mean absolute errors, RMSE, and R2 scores for each model
                                                                                                                                                                                        9 G
        evaluation_scores = []
[12]: # Define the models to evaluate
        models = [
            rets = [
('Linear Regression', LinearRegression()),
('Random Forest', RandomForestRegressor(random_state=35)),
('Gradient Boosting', GradientBoostingRegressor(random_state=35)),
('Support Vector Regression', SVR())
                                                                                                                                                                                        9 G
           for model_name, model in models:
               fold_mse_scores = []
               fold_mae_scores = []
               fold_rmse_scores = []
               fold_r2_scores = []
                                                                                                                                                                                           9 G
  [14]: # Per
          for train_index, test_index in kf.split(X_encoded):
              X_train, X_test = X_encoded.iloc[train_index], X_encoded.iloc[test_index]
y_train, y_test = y.iloc[train_index], y.iloc[test_index]
                                                                                                                                                                                           Q 1
                                                                                                                                                                                           © ©
          model.fit(X_train, y_train)
  [15]: <sub>SVR</sub> 1 0
          SVR()
.6]: # Make predictions on the testing data
                                                                                                                                                                                      ( G
     y_pred = model.predict(X_test)
7]: # Calculate the evaluation scores for this fold
     fold_mse = mean_squared_error(y_test, y_pred)
      fold_mae = mean_absolute_error(y_test, y_pred)
     fold_rmse = root_mean_squared_error(y_test, y_pred) # Use the new function
     fold_r2 = r2_score(y_test, y_pred)
                                                                                                                                                                                     9 G
8]:
     fold_mse_scores.append(fold_mse)
      fold_mae_scores.append(fold_mae)
     fold_rmse_scores.append(fold_rmse)
fold_r2_scores.append(fold_r2)
                                                                                                                                                                                      © G
9]:
     avg_mse = sum(fold_mse_scores) / len(fold_mse_scores)
avg_mae = sum(fold_mae_scores) / len(fold_mae_scores)
     avg_rmse = sum(fold_rmse_scores) / len(fold_rmse_scores)
avg_r2 = sum(fold_r2_scores) / len(fold_r2_scores)
                                                                                                                                                                                      ♀ ⓒ
```

```
avg_n2 = sum(fold_n2_scores) / len(fold_n2_scores)

[20]: # Append the average evaluation scores to the list
evaluation_scores_append((model_name, avg_mse, avg_mse, avg_rmse, avg_r2))

[21]: # Sort the models based on their average mean squared error in ascending order
evaluation_scores.sort(key=lambda x: x[1])

[22]: # Get the best performing model and its corresponding evaluation scores
best_model, best_mse, best_mae, best_r2 = evaluation_scores[0]

[24]: # Print the best_performing model and its evaluation scores
print(f"Best Model: (best_model)")
print(f"Best Model: (best_model)")
print(f"Best Model: (best_mea: 2f)")
print(f"Root Mean Squared Error: (best_mse: 2f)")
print(f"Root Mean Squared Error: (best_mse: 2f)")

Best Model: Support Vector Regression
Mean Squared Error: 7081881.56
Mean Absolute Error: 7081881.56
Mean Absolute Error: 2661.14
R2 Score: 0.00%
```

THE RESULT FOR merging solution: