

# **Data Science III with python (Class notes)**

**STAT 303-3**

Lizhen Shi

2025-04-02

# Table of contents

<b>Preface</b>	<b>8</b>
<b>I Bias &amp; Variance; KNN</b>	<b>9</b>
<b>1 Bias-variance tradeoff</b>	<b>10</b>
1.1 Simple model (Less flexible) . . . . .	10
1.2 Complex model (more flexible) . . . . .	13
<b>2 KNN</b>	<b>16</b>
2.1 KNN for regression . . . . .	16
2.2 Feature Scaling in KNN . . . . .	18
2.3 Hyperparameters in KNN . . . . .	19
2.3.1 Tuning $k$ in KNN . . . . .	20
2.3.2 Tuning Other KNN Hyperparameters . . . . .	23
2.4 Hyperparameter Tuning . . . . .	27
<b>3 Hyperparameter tuning</b>	<b>34</b>
3.1 GridSearchCV . . . . .	35
3.2 RandomizedSearchCV() . . . . .	38
3.3 BayesSearchCV() . . . . .	40
3.3.1 Diagonosis of cross-validated score optimization . . . . .	43
3.3.2 Live monitoring of cross-validated score . . . . .	48
3.4 cross_validate() . . . . .	49
<b>II Tree based models</b>	<b>53</b>
<b>4 Regression trees</b>	<b>54</b>
4.1 Native Support for Missing Values . . . . .	55
4.1.1 Build a regression tree using mileage as the solo predictor . . . . .	56
4.2 Building regression trees . . . . .	57
4.2.1 Using only mileage feature . . . . .	57
4.2.2 Using mileage and brand as predictors . . . . .	59
4.2.3 Using all predictors . . . . .	61

4.3	Key Hyperparameters in Decision Tree . . . . .	63
4.3.1	Underfitting . . . . .	63
4.3.2	Overfitting . . . . .	63
4.3.3	<code>max_depth</code> . . . . .	63
4.3.4	<code>min_samples_split</code> . . . . .	63
4.3.5	<code>min_samples_leaf</code> . . . . .	64
4.3.6	<code>max_features</code> . . . . .	64
4.3.7	Output feature importance . . . . .	67
4.4	Cost-Complexity Pruning ( <code>ccp_alpha</code> ) . . . . .	69
4.4.1	Key Idea . . . . .	69
4.4.2	Parameter: <code>ccp_alpha</code> in scikit-learn . . . . .	69
<b>5</b>	<b>Classification trees</b>	<b>74</b>
5.1	Building a Classification Tree . . . . .	75
5.2	Pre-pruning: Hyperparameters Tuning . . . . .	77
5.2.1	Gini or entropy . . . . .	79
5.3	Post-pruning: Cost complexity pruning . . . . .	81
5.3.1	step 1: calculate the cost complexity pruning path . . . . .	82
5.3.2	step 2: Create trees with different <code>ccp_alpha</code> values and evaluate their performance . . . . .	82
5.3.3	Step 3: Visualize the results . . . . .	82
5.3.4	Step 4: Create the final model with the optimal alpha . . . . .	85
5.4	Feature Importance in Decision Trees . . . . .	86
5.4.1	Do We Need Feature Selection or Regularization with Tree Models? . .	88
5.4.2	Bottom Line . . . . .	89
5.5	Next Lecture . . . . .	89
<b>6</b>	<b>Bagging</b>	<b>90</b>
6.1	Bagging: A Variance Reduction Technique . . . . .	90
6.2	Bagging Regression Trees . . . . .	91
6.3	Bagging Doesn't Reduce Bias . . . . .	95
6.4	Model Performance vs. Number of Trees . . . . .	97
6.5	OOB Sample and OOB Score in Bagging . . . . .	99
6.5.1	What is an OOB Sample? . . . . .	99
6.5.2	What is OOB Score? . . . . .	100
6.6	Bagging Hyperparameter Tuning . . . . .	102
6.6.1	Tuning with Cross-Validation . . . . .	103
6.6.2	Tuning with Out-of-Bag (OOB) Score . . . . .	104
6.6.3	Comparing Hyperparameter Tuning: Cross-Validation vs. OOB Score .	107
6.7	Bagging Classification Trees . . . . .	108

<b>7</b>	<b>Random Forests</b>	<b>110</b>
7.1	Motivation: Bagging Revisited . . . . .	110
7.1.1	Let's build a single decision tree and output its performance . . . . .	111
7.1.2	Let's Build a Bagging Tree with Bootstrap Sampling to Reduce Variance	112
7.1.3	Let's Build a Bagging Tree Without Bootstrap Sampling . . . . .	113
7.1.4	Why Does Bagging Without Bootstrap Perform Worse? . . . . .	115
7.1.5	Why Can Bagging Without Bootstrap Still Show Slight Improvement?	115
7.2	Random Forest . . . . .	115
7.2.1	Idea Behind Random Forest . . . . .	116
7.2.2	Key Hyperparameter Comparison . . . . .	116
7.2.3	Let's Build a Random Forest Model Using the Default Settings . . . . .	116
7.2.4	Let's Build a Random Forest Model with <code>sqrt</code> <code>max_features</code> . . . . .	117
7.3	Let's Explore How <code>max_features</code> Affects Performance . . . . .	119
7.4	Other Hyperparameters in Random Forest . . . . .	121
7.4.1	Why Bagging Uses Unpruned Trees . . . . .	121
7.4.2	Hyperparameters That Control Tree Complexity in Random Forest . . .	121
7.4.3	Why Random Forest Often Limits Tree Depth . . . . .	122
7.4.4	Let's Tune Multiple Hyperparameters Simultaneously Using Cross-Validation . . . . .	122
7.5	Feature Importance in Random Forest . . . . .	124
7.5.1	How Feature Importance Is Calculated . . . . .	124
7.5.2	Accessing Feature Importances . . . . .	124
7.6	In Summary . . . . .	125
7.7	Next Step . . . . .	126
<b>8</b>	<b>Adaptive Boosting</b>	<b>127</b>
8.1	What is AdaBoost? . . . . .	128
8.2	AdaBoost Intuition . . . . .	128
8.3	How AdaBoost Works (High-Level Steps) . . . . .	128
8.4	Key Hyperparameters in AdaBoost . . . . .	128
8.5	AdaBoost for Regression . . . . .	129
8.5.1	Let's build a <code>adaboost</code> regressor model with default settings . . . . .	130
8.5.2	Impact of Tree Depth . . . . .	132
8.5.3	Impact of Learning Rate . . . . .	134
8.5.4	Impact of Number of Trees in Boosting . . . . .	137
8.5.5	Tuning Hyperparameters Simultaneously . . . . .	140
8.5.6	Using Optuna for Hyperparameter Tuning . . . . .	149
<b>9</b>	<b>Gradient Boosting</b>	<b>154</b>
9.1	Core Hyperparameters in Gradient Boosting . . . . .	154
9.1.1	Effect of Number of Trees on Cross-Validation Error . . . . .	157
9.1.2	Early stopping in Gradient Boosting . . . . .	158
9.1.3	Effect of Learning Rate on Cross-Validation Error . . . . .	162

9.1.4	Learning Rate and Number of Trees Are Closely Linked . . . . .	164
9.1.5	Effect of Stochastic Gradient Boosting on Cross-Validation Error . . . .	165
9.1.6	Effect of Tree Complexity on Cross-Validation Error (Not Tuned Here)	167
9.1.7	Loss Function ( <code>loss</code> ) . . . . .	168
9.2	Joint Hyperparameter Tuning in Gradient Boosting for Regression . . . . .	168
9.2.1	Using <code>BayesSearchCV</code> . . . . .	168
9.2.2	Tuning with Optuna . . . . .	173
9.3	Scaling Gradient Boosting for Large Datasets and Smarter Tuning . . . . .	176
9.4	Independent Study . . . . .	177
<b>10</b>	<b>XGBoost</b>	<b>178</b>
10.1	Hyperparameters . . . . .	178
10.2	XGBoost for regression . . . . .	180
10.2.1	Number of trees vs cross validation error . . . . .	180
10.2.2	Depth of tree vs cross validation error . . . . .	181
10.2.3	Learning rate vs cross validation error . . . . .	183
10.2.4	Regularization ( <code>reg_lambda</code> ) vs cross validation error . . . . .	185
10.2.5	Regularization ( <code>gamma</code> ) vs cross validation error . . . . .	188
10.2.6	Tuning XGboost regressor . . . . .	190
10.2.7	Early stopping with XGBoost . . . . .	193
10.3	XGBoost for classification . . . . .	195
10.3.1	Precision & recall vs <code>scale_pos_weight</code> . . . . .	200
<b>11</b>	<b>LightGBM and CatBoost</b>	<b>203</b>
11.1	LightGBM . . . . .	204
11.1.1	LightGBM for regression . . . . .	205
11.1.2	LightGBM vs XGBoost . . . . .	208
11.2	CatBoost . . . . .	209
11.2.1	CatBoost for regression . . . . .	210
11.2.2	Target encoding with CatBoost . . . . .	211
11.2.3	CatBoost vs XGBoost . . . . .	212
11.2.4	Tuning <code>CatBoostRegressor</code> . . . . .	213
11.2.5	Tuning Tips . . . . .	215
<b>12</b>	<b>Ensemble modeling</b>	<b>216</b>
12.1	Ensembling regression models . . . . .	218
12.1.1	Voting Regressor . . . . .	218
12.1.2	Stacking Regressor . . . . .	221
12.2	Ensembling classification models . . . . .	225
	AdaBoost . . . . .	226
	Gradient Boosting . . . . .	226
	XGBoost . . . . .	227
12.2.1	Voting classifier - hard voting . . . . .	228

12.2.2	Voting classifier - soft voting . . . . .	229
12.2.3	Stacking classifier . . . . .	229
12.2.4	Tuning all models simultaneously . . . . .	231
12.3	Ensembling models based on different sets of predictors . . . . .	232

## Appendices 236

### A Assignment 1 236

Instructions . . . . .	236
A.1 1) Bias-Variance Trade-off for Regression ( <b>50 points</b> ) . . . . .	236
A.1.1 a) Define the True Relationship (Signal) . . . . .	237
A.1.2 b) Generate Test Set (No Noise) . . . . .	237
A.1.3 c) Initialize Results DataFrame . . . . .	238
A.1.4 d) Generate Training Sets (With Noise) . . . . .	238
A.1.5 e) Visualize Bias-Variance Decomposition . . . . .	239
A.1.6 f) Identify the Optimal Model . . . . .	240
A.2 2) Building a Low-Bias, Low-Variance Model via Regularization (50 points) . .	240
A.2.1 a) Why Regularization? . . . . .	240
A.2.2 b) Which Degrees to Exclude? . . . . .	241
A.2.3 c) Apply Ridge Regularization . . . . .	241
A.2.4 d) Visualize Regularized Results . . . . .	241
A.2.5 e) Evaluate the Regularized Model . . . . .	242
A.2.6 f) Interpreting the Impact of Regularization . . . . .	242

### B Assignment 2 243

Instructions . . . . .	243
B.1 Optimizing KNN for Classification (71 points) . . . . .	244
B.1.1 a) Load the Dataset ( <b>1 point</b> ) . . . . .	244
B.1.2 b) Define Predictor and Response Variables ( <b>1 point</b> ) . . . . .	244
B.1.3 c) Split the Data into Training and Test Sets ( <b>1 points</b> ) . . . . .	244
B.1.4 d) Check Class Ratios ( <b>2 points</b> ) . . . . .	245
B.1.5 e) Scale the Dataset ( <b>2 points</b> ) . . . . .	245
B.1.6 f) Set Up Cross-Validation ( <b>2 points</b> ) . . . . .	245
B.1.7 g) Tune K for KNN Using Cross-Validation ( <b>12 points</b> ) . . . . .	245
B.1.8 h) Optimize the Classification Threshold ( <b>4 points</b> ) . . . . .	246
B.1.9 i) Evaluate the Tuning Method ( <b>2 points</b> ) . . . . .	247
B.1.10 j) Evaluate Tuned Classifier on Test Set ( <b>3 points</b> ) . . . . .	247
B.1.11 k) Jointly Tune K and Threshold ( <b>6 points</b> ) . . . . .	247
B.1.12 l) Visualize Cross-Validation Results with a Heatmap ( <b>3 points</b> ) . . . .	247
B.1.13 m) Compare Joint vs. Sequential Tuning Results ( <b>4 points</b> ) . . . . .	248
B.1.14 n) Evaluate Final Tuned Model on Test Set ( <b>3 points</b> ) . . . . .	248
B.1.15 o) Compare Tuning Strategies and Computational Cost ( <b>3 points</b> ) . .	248

B.1.16	p) Tune K Using Multiple Metrics ( <b>5 points</b> ) . . . . .	248
B.1.17	q) Optimize for Recall with Precision Constraint ( <b>4 point</b> ) . . . . .	249
B.1.18	r) Tune Threshold for Maximum Recall ( <b>3 point</b> ) . . . . .	249
B.1.19	s) Evaluate Precision-Optimized Model on Test Set ( <b>2 points</b> ) . . . . .	249
B.1.20	t) Final Reflection: Comparing Tuning Strategies ( <b>3 points</b> ) . . . . .	250
B.2	Tuning a KNN Regressor on Bank Loan Data (32 points) . . . . .	250
B.2.1	a) Split, Scale, and Tune a KNN Regressor ( <b>15 point</b> ) . . . . .	251
B.2.2	b) Compare Tuning Approaches ( <b>1 point</b> ) . . . . .	253
B.2.3	c) Feature Selection and Hyperparameter Tuning with GridSearchCV ( <b>15 point</b> ) . . . . .	253
B.2.4	d) Compare Feature Selection Approaches ( <b>1 point</b> ) . . . . .	254
<b>C</b>	<b>Assignment 3</b>	<b>255</b>
	Instructions . . . . .	255
C.1	Regression Problem - Miami housing . . . . .	255
C.1.1	a) Data preparation . . . . .	255
C.1.2	b) Baseline Decision Tree Model . . . . .	256
C.1.3	c) Tune the Decision Tree Model . . . . .	256
C.1.4	d) Bagged Decision Trees with Out-of-Bag Evaluation . . . . .	256
C.1.5	e) Bagged Decision Trees Without Bootstrapping . . . . .	257
C.1.6	f) Bagged Decision Trees with Feature Bootstrapping Only . . . . .	257
C.1.7	g) Tuning a Bagged Tree Model . . . . .	258
C.1.8	h) Random Forest . . . . .	258
C.2	Classification - Term deposit . . . . .	260
C.2.1	a) Data Preparation . . . . .	260
C.2.2	b) Random Forest for Term Deposit Subscription Prediction . . . . .	261
C.3	Predictor Transformations in Trees . . . . .	262
<b>D</b>	<b>Datasets, assignment and project files</b>	<b>263</b>

# Preface

This book serves as the course notes for STAT 303 Sec20, Spring 2025 at Northwestern University. To enhance your understanding of the material, you are expected to read the textbook before using these notes.

It is an evolving resource designed to support the course's learning objectives. This edition builds upon the foundational work of Professor Arvind Krishna, whose contributions have provided a strong framework for this resource. We are deeply grateful for his efforts, which continue to shape the book's development.

Throughout the quarter, the content will be updated and refined in real time to enhance clarity, depth, and relevance. These modifications ensure alignment with current teaching objectives and methodologies.

As a living document, this book welcomes feedback, suggestions, and contributions from students, instructors, and the broader academic community. Your input helps improve its quality and effectiveness.

Thank you for being part of this journey—we hope this resource serves as a valuable guide in your learning.



## **Part I**

# **Bias & Variance; KNN**

# 1 Bias-variance tradeoff

*Read section 2.2.2 of the book before using these notes.*

*Note that in this course, lecture notes are not sufficient, you must read the book for better understanding. Lecture notes are just implementing the concepts of the book on a dataset, but not explaining the concepts elaborately.*

In this chapter, we will show that a flexible model is likely to have high variance and low bias, while a relatively less flexible model is likely to have a high bias and low variance.

The examples considered below are motivated from the examples shown in the documentation of the `bias_variance_decomp()` function from the `mlxtend` library. We will first manually compute the bias and variance for understanding of the concept. Later, we will show application of the `bias_variance_decomp()` function to estimate bias and variance.

```
# Importing necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.tree import DecisionTreeRegressor
sns.set(font_scale=1.35)
```

## 1.1 Simple model (Less flexible)

Let us consider a linear regression model as the less-flexible (*or relatively simple*) model.

We will first simulate the test dataset for which we will compute the bias and variance.

```
np.random.seed(101)

# Simulating predictor values of test data
xtest = np.random.uniform(-15, 10, 200)
```

```

# Assuming the true mean response is square of the predictor value
fxtest = xtest**2

# Simulating noiseless test response
ytest = fxtest

# We will find bias and variance using a linear regression model for prediction
model = LinearRegression()

# Visualizing the data and the true mean response
sns.scatterplot(x = xtest, y = ytest)
sns.lineplot(x = xtest, y = fxtest, color = 'grey', linewidth = 2)

# Initializing objects to store predictions and mean squared error
# of 100 models developed on 100 distinct training datasets samples
pred_test = []; mse_test = []

# Iterating over each of the 100 models
for i in range(100):
    np.random.seed(i)

    # Simulating the ith training data
    x = np.random.uniform(-15, 10, 200)
    fx = x**2
    y = fx + np.random.normal(0, 10, 200)

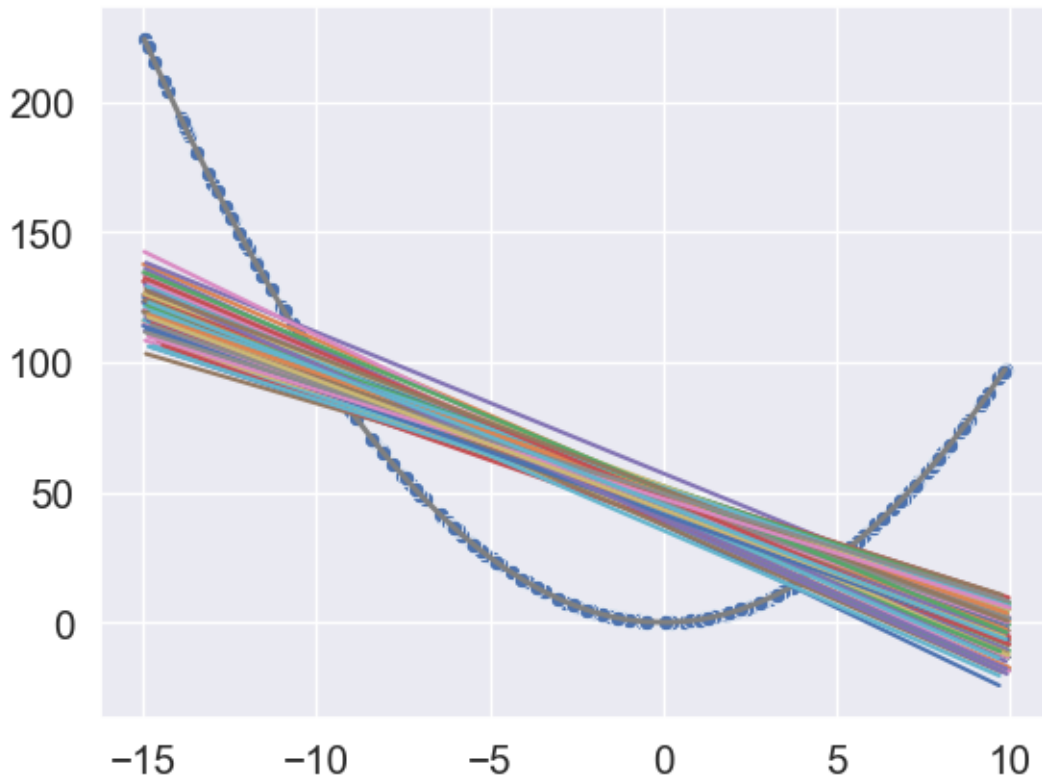
    # Fitting the ith model on the ith training data
    model.fit(x.reshape(-1,1), y)

    # Plotting the ith model
    sns.lineplot(x = x, y = model.predict(x.reshape(-1,1)))

    # Storing the predictions of the ith model on test data
    pred_test.append(model.predict(xtest.reshape(-1,1)))

    # Storing the mean squared error of the ith model on test data
    mse_test.append(mean_squared_error(model.predict(xtest.reshape(-1,1)), ytest))

```



The above plots show that the 100 models seem to have low variance, but high bias. Note that the bias is low only around a couple of points ( $x = -10$  &  $x = 5$ ).

Let us compute the average squared bias over all the test data points.

```
mean_pred = np.array(pred_test).mean(axis = 0)
sq_bias = ((mean_pred - fxtest)**2).mean()
sq_bias
```

2042.104126728109

Let us compute the average variance over all the test data points.

```
mean_var = np.array(pred_test).var(axis = 0).mean()
mean_var
```

28.37397844429763

Let us compute the mean squared error over all the test data points.

```
np.array(mse_test).mean()
```

2070.4781051724062

Note that the mean squared error should be the same as the sum of squared bias and variance

The sum of squared bias and model variance is:

```
sq_bias + mean_var
```

2070.4781051724067

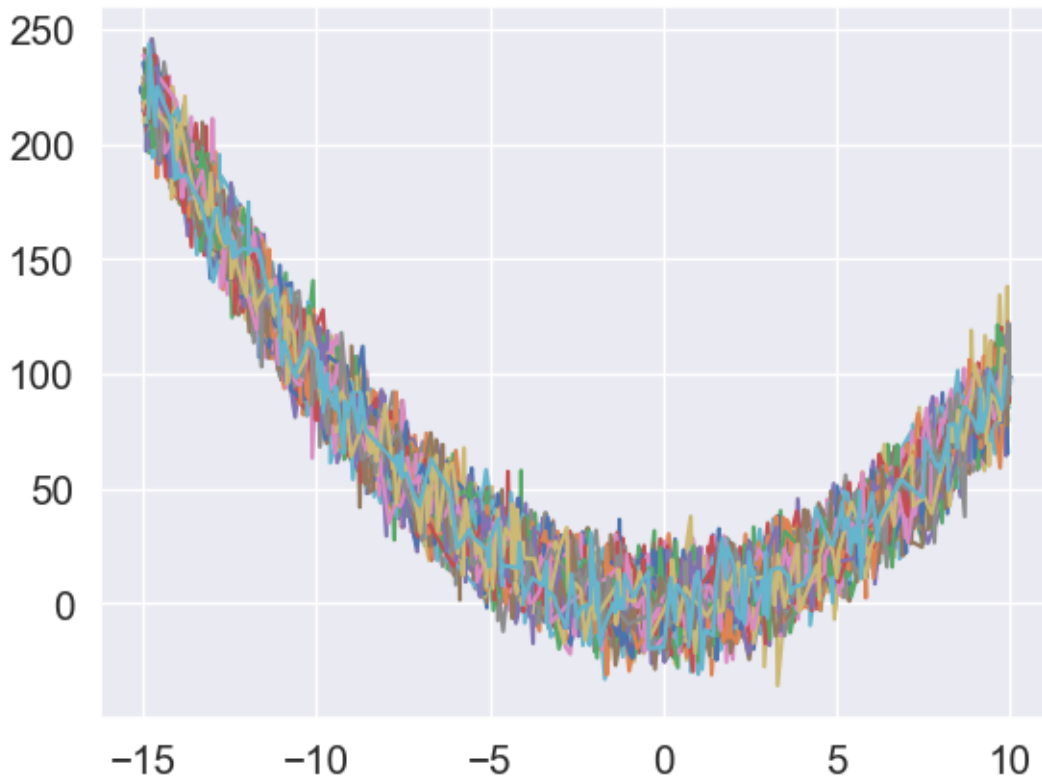
Note that this is exactly the same as the mean squared error computed above as we are developing a finite number of models, and making predictions on a finite number of test data points.

## 1.2 Complex model (more flexible)

Let us consider a decision tree as the more flexible model.

```
np.random.seed(101)
xtest = np.random.uniform(-15, 10, 200)
fxtest = xtest**2
ytest = fxtest
model = DecisionTreeRegressor()
```

```
sns.scatterplot(x = xtest, y = ytest)
sns.lineplot(x = xtest, y = fxtest, color = 'grey', linewidth = 2)
pred_test = []; mse_test = []
for i in range(100):
    np.random.seed(i)
    x = np.random.uniform(-15, 10, 200)
    fx = x**2
    y = fx + np.random.normal(0, 10, 200)
    model.fit(x.reshape(-1,1), y)
    sns.lineplot(x = x, y = model.predict(x.reshape(-1,1)))
    pred_test.append(model.predict(xtest.reshape(-1,1)))
    mse_test.append(mean_squared_error(model.predict(xtest.reshape(-1,1)), ytest))
```



The above plots show that the 100 models seem to have high variance, but low bias.

Let us compute the average squared bias over all the test data points.

```
mean_pred = np.array(pred_test).mean(axis = 0)
sq_bias = ((mean_pred - fxtest)**2).mean()
sq_bias
```

1.3117561629333938

Let us compute the average model variance over all the test data points.

```
mean_var = np.array(pred_test).var(axis = 0).mean()
mean_var
```

102.5226748977198

Let us compute the average mean squared error over all the test data points.

```
np.array(mse_test).mean()
```

```
103.83443106065317
```

Note that the above error is still the same as the sum of the squared bias, model variance and the irreducible error.

Note that the relatively more flexible model has a higher variance, but lower bias as compared to the less flexible linear model. This will typically be the case, but may not be true in all scenarios. We will discuss one such scenario later.

## 2 KNN

*Read section 4.7.6 of the book before using these notes.*

*Note that in this course, lecture notes are not sufficient, you must read the book for better understanding. Lecture notes are just implementing the concepts of the book on a dataset, but not explaining the concepts elaborately.*

```
# Importing necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(font_scale=1.35)

from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsRegressor, KNeighborsClassifier
from sklearn.model_selection import cross_val_score, GridSearchCV, cross_val_predict, KFold,

from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, FunctionTransformer
from sklearn.metrics import root_mean_squared_error, r2_score
```

### 2.1 KNN for regression

```
# Load the dataset
car = pd.read_csv('Datasets/car.csv')

# Split the dataset into features and target variable
X = car.drop(columns=['price'])
y = car['price']

# split the dataset into training and testing sets
```



```

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# extract the categorical columns and put them in a list
cat_cols = X.select_dtypes(include=['object']).columns.tolist()

# extract the numerical columns and put them in a list
num_cols = X.select_dtypes(include=['int64', 'float64']).columns.tolist()

# First transform categorical variables
preprocessor = ColumnTransformer(
    transformers=[
        ('num', 'passthrough', num_cols), # Just pass numerical features through
        ('cat', OneHotEncoder(handle_unknown='ignore', sparse_output=False), cat_cols)
    ])

# Create pipeline that scales all features together
pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('scaler', StandardScaler()), # Scale everything together
    ('knn', KNeighborsRegressor(n_neighbors=5))
])

# Fit the pipeline to the training data
pipeline.fit(X_train, y_train)
# Predict on the test data
y_pred = pipeline.predict(X_test)
# Calculate RMSE
rmse = root_mean_squared_error(y_test, y_pred)
print(f"RMSE: {rmse:.2f}")
print(f"R2 Score: {pipeline.score(X_test, y_test):.2f}")

```

RMSE: 4364.84

R<sup>2</sup> Score: 0.94

```

# show the features in the numerical transformer
pipeline.named_steps['preprocessor'].transformers_[0][1].get_feature_names_out()
print("numerical features in the pipeline:", pipeline.named_steps['preprocessor'].transformer

# show the features in the categorical transformer
pipeline.named_steps['preprocessor'].transformers_[1][1].get_feature_names_out()
print("categorical features in the pipeline:", pipeline.named_steps['preprocessor'].transformer

```

```

numerical features in the pipeline: ['year' 'mileage' 'tax' 'mpg' 'engineSize']
categorical features in the pipeline: ['brand_audi' 'brand_bmw' 'brand_ford' 'brand_hyundai'
'brand_skoda' 'brand_toyota' 'brand_vauxhall' 'brand_vw'
'model_ 6 Series' 'model_ 7 Series' 'model_ 8 Series' 'model_ A7'
'model_ A8' 'model_ Agila' 'model_ Amarok' 'model_ Antara'
'model_ Arteon' 'model_ Avensis' 'model_ Beetle' 'model_ CC'
'model_ CLA Class' 'model_ CLK' 'model_ CLS Class' 'model_ Caddy'
'model_ Caddy Life' 'model_ Caddy Maxi Life' 'model_ California'
'model_ Camry' 'model_ Caravelle' 'model_ Combo Life' 'model_ Edge'
'model_ Eos' 'model_ Fusion' 'model_ G Class' 'model_ GL Class'
'model_ GLB Class' 'model_ GLS Class' 'model_ GT86' 'model_ GTC'
'model_ Galaxy' 'model_ Getz' 'model_ Grand C-MAX'
'model_ Grand Tourneo Connect' 'model_ Hilux' 'model_ I40' 'model_ I800'
'model_ IQ' 'model_ IX20' 'model_ IX35' 'model_ Jetta' 'model_ KA'
'model_ Kamiq' 'model_ Land Cruiser' 'model_ M Class' 'model_ M2'
'model_ M3' 'model_ M4' 'model_ M5' 'model_ M6' 'model_ Mustang'
'model_ PROACE VERSO' 'model_ Prius' 'model_ Puma' 'model_ Q8'
'model_ R8' 'model_ RS3' 'model_ RS4' 'model_ RS5' 'model_ RS6'
'model_ Rapid' 'model_ Roomster' 'model_ S Class' 'model_ S3' 'model_ S4'
'model_ SLK' 'model_ SQ5' 'model_ SQ7' 'model_ Santa Fe' 'model_ Scala'
'model_ Scirocco' 'model_ Shuttle' 'model_ Supra'
'model_ Tiguan Allspace' 'model_ Tourneo Connect' 'model_ Tourneo Custom'
'model_ V Class' 'model_ Verso' 'model_ Vivaro' 'model_ X-CLASS'
'model_ X4' 'model_ X6' 'model_ X7' 'model_ Yeti' 'model_ Z3' 'model_ Z4'
'model_ Zafira Tourer' 'model_ i3' 'model_ i8' 'transmission_Automatic'
'transmission_Manual' 'transmission_Other' 'transmission_Semi-Auto'
'fuelType_Diesel' 'fuelType_Electric' 'fuelType_Hybrid' 'fuelType_Other'
'fuelType_Petrol']

```

## 2.2 Feature Scaling in KNN

**Feature scaling** is essential when using **K-Nearest Neighbors (KNN)** because the algorithm relies on calculating distances between data points. If features are measured on different scales (e.g., **mileage** in thousands and **mpg** in tens), the features with larger numeric ranges can dominate the distance calculations and distort the results.

To ensure that all features contribute equally, it's important to **standardize or normalize** them before applying KNN. Common scaling techniques include:

- **Standardization** (zero mean, unit variance) using `StandardScaler`

- **Min-max scaling** to bring values into the  $[0, 1]$  range

Without scaling, KNN may produce biased or misleading predictions.

The example below illustrates how the same KNN model performs **without feature scaling**, highlighting the importance of preprocessing your data.

```
preprocessor_no_scaling = ColumnTransformer(
    transformers=[
        ('num', 'passthrough', num_cols), # Pass numerical features through without scaling
        ('cat', OneHotEncoder(handle_unknown='ignore'), cat_cols) # Only one-hot encode cat
    ])

# Create pipeline without any scaling
pipeline_no_scaling = Pipeline(steps=[
    ('preprocessor', preprocessor_no_scaling),
    ('knn', KNeighborsRegressor(n_neighbors=5))
])

# Fit the pipeline
pipeline_no_scaling.fit(X_train, y_train)

# Evaluate
y_pred_no_scaling = pipeline_no_scaling.predict(X_test)

rmse_no_scaling = root_mean_squared_error(y_test, y_pred_no_scaling)
print(f"RMSE without scaling: {rmse_no_scaling:.2f}")
print(f"R2 Score without scaling: {pipeline_no_scaling.score(X_test, y_test):.2f}")
```

RMSE without scaling: 13758.38

R<sup>2</sup> Score without scaling: 0.35

## 2.3 Hyperparameters in KNN

The most important hyperparameter in K-Nearest Neighbors (KNN) is  $k$ , which determines the number of neighbors considered when making predictions. Tuning  $k$  helps balance the model's **bias and variance**:

- A **small  $k$**  (e.g., 1 or 3) can lead to **low bias but high variance**, making the model sensitive to noise in the training data.
- A **large  $k$**  results in **higher bias but lower variance**, producing smoother predictions that may underfit the data.

### 2.3.1 Tuning $k$ in KNN

To find the optimal value of  $k$ , it's common to use **cross-validation**, which evaluates model performance on different subsets of the data. A popular tool for this is **GridSearchCV**, which automates the search process by testing multiple values of  $k$  using cross-validation behind the scenes. It selects the value of  $k$  that minimizes prediction error on unseen data—helping you achieve a good balance between underfitting and overfitting.

```
# Create parameter grid for k values
param_grid = {
    'knn__n_neighbors': list(range(1, 20)) # Test k values from 1 to 20
}

# Set up GridSearchCV
grid_search = GridSearchCV(
    estimator=pipeline,
    param_grid=param_grid,
    cv=5, # 5-fold cross-validation
    scoring='neg_root_mean_squared_error', # Optimize for RMSE
    n_jobs=-1, # Use all available cores
    verbose=1
)

# Fit grid search
print("Tuning k parameter...")
grid_search.fit(X_train, y_train)

# Get best parameters and results
best_k = grid_search.best_params_['knn__n_neighbors']
best_score = -grid_search.best_score_ # Convert back from negative RMSE

print(f"Best k: {best_k}")
print(f"Best CV RMSE: {best_score:.2f}")

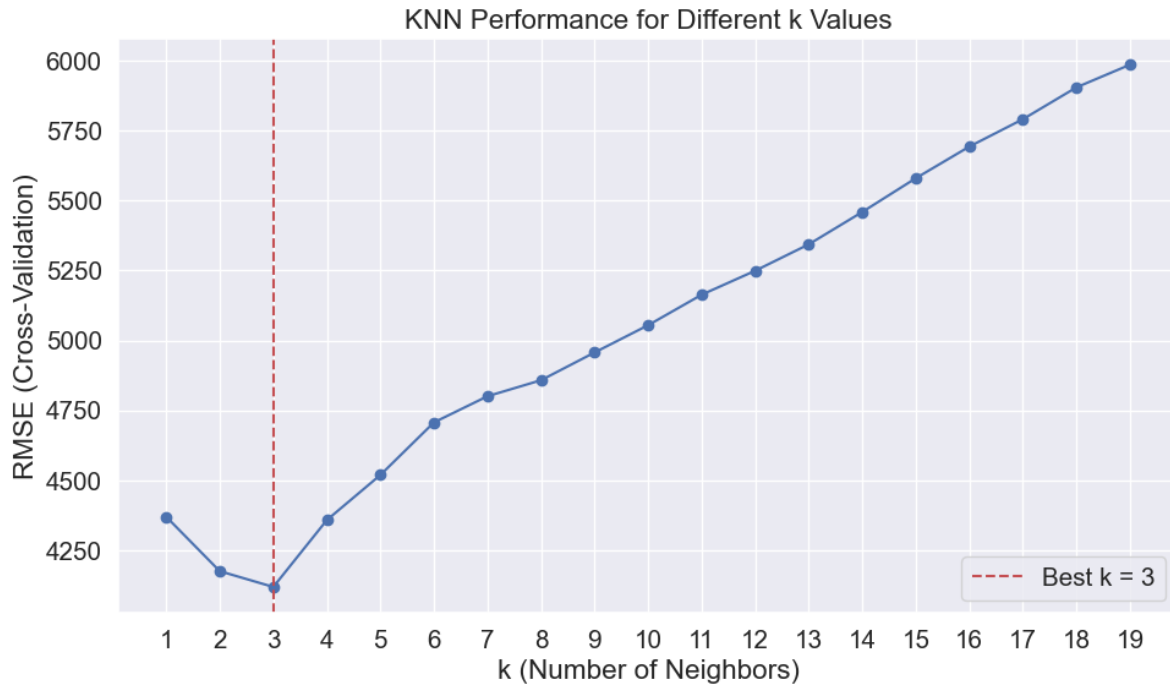
# Evaluate on test set using best model
best_model = grid_search.best_estimator_
y_pred = best_model.predict(X_test)
test_rmse = root_mean_squared_error(y_test, y_pred)
test_r2 = r2_score(y_test, y_pred)

print(f"Test RMSE with k={best_k}: {test_rmse:.2f}")
print(f"Test R2 Score with k={best_k}: {test_r2:.2f}")
```

Tuning k parameter...  
Fitting 5 folds for each of 19 candidates, totalling 95 fits  
Best k: 3  
Best CV RMSE: 4117.42  
Test RMSE with k=3: 4051.06  
Test R<sup>2</sup> Score with k=3: 0.94

```
# Plot performance across different k values
cv_results = grid_search.cv_results_
k_values = param_grid['knn__n_neighbors']
mean_rmse = -cv_results['mean_test_score']

plt.figure(figsize=(10, 6))
plt.plot(k_values, mean_rmse, marker='o')
plt.xlabel('k (Number of Neighbors)')
plt.ylabel('RMSE (Cross-Validation)')
plt.title('KNN Performance for Different k Values')
plt.grid(True)
plt.xticks(k_values)
plt.axvline(x=best_k, color='r', linestyle='--', label=f'Best k = {best_k}')
plt.legend()
plt.tight_layout()
plt.show()
```



The distances and the indices of the nearest K observations to each test observation can be obtained using the `kneighbors()` method.

```
# Get the KNN estimator from the pipeline
knn_estimator = best_model.named_steps['knn']

# Get indices of K-nearest neighbors for each test observation
neighbor_indices = knn_estimator.kneighbors(best_model.named_steps['preprocessor'].transform(test_data),
                                           return_distance=False)

# neighbor_indices will contain the indices of the K nearest neighbors for each test observation
# Note: The indices are relative to the training set, not the test set.
# To get the actual neighbor observations, you can use these indices to index into the training set.
# For example, to get the actual neighbor observations for the first test observation:
neighbors = X_train.iloc[neighbor_indices[0]]
neighbors
```

	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize
4580	merc	V Class	2010	Automatic	259000	Diesel	540	30.8345	3.0
5651	merc	CLK	2003	Automatic	185000	Petrol	330	18.0803	4.3
3961	vw	Caravelle	2006	Manual	178000	Diesel	325	34.5738	2.5

### 2.3.2 Tuning Other KNN Hyperparameters

In addition to the number of neighbors ( $k$ ), KNN has several other important hyperparameters that can significantly affect the model's performance. Fine-tuning these settings helps you get the most out of the algorithm. Key hyperparameters include:

- **weights**: Determines how the neighbors contribute to the prediction.
  - 'uniform': All neighbors are weighted equally (default).
  - 'distance': Closer neighbors have more influence.
  - Choosing 'distance' can improve performance, especially when data points are unevenly distributed.
- **metric**: Defines the distance function used to measure similarity between data points.
  - 'minkowski' (default) is a general-purpose metric that includes both Euclidean and Manhattan distances.
  - Other options include 'euclidean', 'manhattan', or even custom distance functions.
- **p**: Used when `metric='minkowski'`.
  - $p=2$  gives **Euclidean distance** (standard for continuous features).
  - $p=1$  gives **Manhattan distance** (useful when features are sparse or grid-based).
- **algorithm**: Controls the method used to compute nearest neighbors.
  - 'auto', 'ball\_tree', 'kd\_tree', or 'brute'.
  - Most users can leave this as 'auto', which lets scikit-learn choose the best algorithm based on the data.

These hyperparameters can be tuned using `GridSearchCV` to find the combination that yields the best performance on validation data.

The model hyperparameters can be obtained using the `get_params()` method. Note that there are other hyperparameters to tune in addition to number of neighbors. However, the number of neighbours may be the most influential hyperparameter in most cases.

```
# Get the best model parameters
best_model.get_params()
```

```

{'memory': None,
 'steps': [('preprocessor',
           ColumnTransformer(transformers=[('num', 'passthrough',
                                           ['year', 'mileage', 'tax', 'mpg',
                                            'engineSize']),
                                           ('cat',
                                            OneHotEncoder(handle_unknown='ignore',
                                                            sparse_output=False),
                                            ['brand', 'model', 'transmission',
                                             'fuelType'])])),
           ('scaler', StandardScaler()),
           ('knn', KNeighborsRegressor(n_neighbors=3))],
 'transform_input': None,
 'verbose': False,
 'preprocessor': ColumnTransformer(transformers=[('num', 'passthrough',
                                                 ['year', 'mileage', 'tax', 'mpg',
                                                  'engineSize']),
                                                 ('cat',
                                                  OneHotEncoder(handle_unknown='ignore',
                                                                    sparse_output=False),
                                                  ['brand', 'model', 'transmission',
                                                   'fuelType'])])),
 'scaler': StandardScaler(),
 'knn': KNeighborsRegressor(n_neighbors=3),
 'preprocessor__force_int_remainder_cols': True,
 'preprocessor__n_jobs': None,
 'preprocessor__remainder': 'drop',
 'preprocessor__sparse_threshold': 0.3,
 'preprocessor__transformer_weights': None,
 'preprocessor__transformers': [('num',
                                'passthrough',
                                ['year', 'mileage', 'tax', 'mpg', 'engineSize']),
                                ('cat',
                                 OneHotEncoder(handle_unknown='ignore', sparse_output=False),
                                 ['brand', 'model', 'transmission', 'fuelType'])],
 'preprocessor__verbose': False,
 'preprocessor__verbose_feature_names_out': True,
 'preprocessor__num': 'passthrough',
 'preprocessor__cat': OneHotEncoder(handle_unknown='ignore', sparse_output=False),
 'preprocessor__cat__categories': 'auto',
 'preprocessor__cat__drop': None,
 'preprocessor__cat__dtype': numpy.float64,
 'preprocessor__cat__feature_name_combiner': 'concat',

```



```

'preprocessor__cat__handle_unknown': 'ignore',
'preprocessor__cat__max_categories': None,
'preprocessor__cat__min_frequency': None,
'preprocessor__cat__sparse_output': False,
'scaler__copy': True,
'scaler__with_mean': True,
'scaler__with_std': True,
'knn__algorithm': 'auto',
'knn__leaf_size': 30,
'knn__metric': 'minkowski',
'knn__metric_params': None,
'knn__n_jobs': None,
'knn__n_neighbors': 3,
'knn__p': 2,
'knn__weights': 'uniform'}

```

```

# Extended parameter grid
param_grid = {
    'knn__n_neighbors': list(range(1, 20, 2)), # Test odd k values from 1 to 19 (step=2 for
    'knn__weights': ['uniform', 'distance'], # Uniform: equal weight; Distance: closer neigh
    'knn__metric': ['euclidean', 'manhattan', 'minkowski'], # Common distance metrics
    'knn__p': [1, 2] # p=1 (Manhattan), p=2 (Euclidean) - only relevant for Minkowski
}

# Set up GridSearchCV
grid_search = GridSearchCV(
    estimator=pipeline,
    param_grid=param_grid,
    cv=5, # 5-fold cross-validation
    scoring='neg_root_mean_squared_error', # Optimize for RMSE
    n_jobs=-1, # Use all available cores
    verbose=1
)

# Fit grid search
print("Tuning KNN hyperparameters...")
grid_search.fit(X_train, y_train)

# Get best parameters and results
best_params = grid_search.best_params_
best_score = -grid_search.best_score_ # Convert negative RMSE to positive

```

```
# Display results
print("\nBest Parameters:")
for param, value in best_params.items():
    print(f"{param}: {value}")
print(f"Best CV RMSE: {best_score:.2f}")

# Evaluate on test set using best model
best_model = grid_search.best_estimator_
y_pred = best_model.predict(X_test)
test_rmse = root_mean_squared_error(y_test, y_pred) # Calculate RMSE
print(f"Test RMSE: {test_rmse:.2f}")
```

Tuning KNN hyperparameters...

Fitting 5 folds for each of 120 candidates, totalling 600 fits

```
Best Parameters:
knn__metric: euclidean
knn__n_neighbors: 3
knn__p: 1
knn__weights: distance
Best CV RMSE: 4001.34
Test RMSE: 3826.94
```

The results for each cross-validation are stored in the `cv_results_` attribute.

```
pd.DataFrame(grid_search.cv_results_).head()
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_knn__metric	param_knn__n_neighbors
0	0.033124	0.003042	0.127347	0.023598	euclidean	1
1	0.035615	0.010407	0.179835	0.013115	euclidean	1
2	0.027877	0.002536	0.148597	0.018612	euclidean	1
3	0.043631	0.016927	0.168392	0.027444	euclidean	1
4	0.043071	0.009615	0.184532	0.042681	euclidean	3

These results can be useful to see if other hyperparameter values are equally good.

```
pd.DataFrame(grid_search.cv_results_).sort_values(by = 'rank_test_score').head()
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_knn__metric	param_k
87	0.038193	0.010690	0.149261	0.050225	minkowski	3
5	0.047902	0.013181	0.185623	0.049865	euclidean	3
7	0.040595	0.005817	0.132290	0.009807	euclidean	3
51	0.034996	0.001900	0.744842	0.052065	manhattan	5
49	0.031465	0.004676	0.718503	0.057517	manhattan	5

The results show that the next two best hyperparameter values yield the same performance as the printed one

## 2.4 Hyperparameter Tuning

We used `GridSearchCV` to tune the hyperparameters of our KNN model above. Given a relatively simple set of hyperparameters and a limited number of combinations, this approach was sufficient to reduce the RMSE.

However, when the number of possible hyperparameter values grows large, `GridSearchCV` can become computationally expensive. In such cases, `RandomizedSearchCV` provides a more efficient alternative by sampling a fixed number of random combinations from the specified hyperparameter space. This makes it well-suited for scenarios with limited computational resources.

### 2.4.0.1 RandomizedSearchCV

```
from sklearn.model_selection import RandomizedSearchCV
from scipy.stats import randint
# Set up RandomizedSearchCV

# Define parameter distributions for randomized search
param_distributions = {
    'knn__n_neighbors': randint(1, 20), # Random ints from 1 to 19
    'knn__weights': ['uniform', 'distance'],
    'knn__metric': ['euclidean', 'manhattan', 'minkowski'],
    'knn__p': [1, 2] # Only relevant for Minkowski
}
```

```
# Set up RandomizedSearchCV
random_search = RandomizedSearchCV(
    estimator=pipeline,
    param_distributions=param_distributions,
    n_iter=30, # Number of random combinations to try
    cv=5,
    scoring='neg_root_mean_squared_error',
    n_jobs=-1,
    random_state=42,
    verbose=1
)
```

```
# Fit randomized search
print("Tuning KNN hyperparameters with RandomizedSearchCV...")
random_search.fit(X_train, y_train)

# Best results
best_params = random_search.best_params_
best_score = -random_search.best_score_

# Display results
print("\nBest Parameters (RandomizedSearchCV):")
for param, value in best_params.items():
    print(f"{param}: {value}")
print(f"Best CV RMSE: {best_score:.2f}")
```

Tuning KNN hyperparameters with RandomizedSearchCV...  
 Fitting 5 folds for each of 30 candidates, totalling 150 fits

```
Best Parameters (RandomizedSearchCV):
knn__metric: manhattan
knn__n_neighbors: 8
knn__p: 2
knn__weights: distance
Best CV RMSE: 4005.70
```

```
# Evaluate on test set
best_model = random_search.best_estimator_
y_pred = best_model.predict(X_test)

# Calculate RMSE
```

```
test_rmse = root_mean_squared_error(y_test, y_pred)
print(f"Test RMSE: {test_rmse:.2f}")
```

Test RMSE: 3811.89

### Why might `RandomizedSearchCV` outperform `GridSearchCV`?

Although `GridSearchCV` systematically evaluates all combinations of hyperparameter values from a predefined grid, it doesn't guarantee the best performance. In some cases, `RandomizedSearchCV` can actually perform better. Here's why:

- **Limited Grid Resolution:**

`GridSearchCV` evaluates only the specific values you include in the grid. If the true optimal value lies between grid points, it may be missed entirely.

- **Broader Exploration:**

`RandomizedSearchCV` samples from distributions (e.g., continuous or discrete ranges), allowing it to explore a wider range of hyperparameter values, including combinations not explicitly considered in a grid.

In this case,

- `list(range(1, 20, 2))` in `GridSearchCV`
- But in `RandomizedSearchCV`, it samples from `randint(1, 20)`

The best `n_neighbors` happens to be 11, only `RandomizedSearchCV` can find it unless you explicitly included it in your grid.

- **Efficiency in High Dimensions:**

In high-dimensional search spaces, the number of combinations in a grid grows exponentially. `RandomizedSearchCV` remains efficient by sampling a fixed number of combinations, avoiding the “curse of dimensionality.”

- **Better Use of Time Budget:**

Given the same computational budget, `RandomizedSearchCV` may cover more diverse regions of the search space and stumble upon better-performing configurations.

In summary, `RandomizedSearchCV` is not only faster but can also lead to better models—especially when the hyperparameter space is large, continuous, or contains irrelevant parameters.

### 2.4.0.2 BayesSearchCV

In addition to these methods, **BayesSearchCV**, based on Bayesian optimization, provides a more intelligent approach to hyperparameter tuning. It models the performance landscape and selects hyperparameter combinations to evaluate based on past results, often requiring fewer evaluations to find optimal or near-optimal values. This makes **BayesSearchCV** a powerful option, especially when training models is costly.

```
# Step 1: Install scikit-optimize if not already installed
!pip install scikit-optimize
```

Collecting scikit-optimize

Downloading scikit\_optimize-0.10.2-py2.py3-none-any.whl.metadata (9.7 kB)

Requirement already satisfied: joblib>=0.11 in c:\users\lsi8012\appdata\local\anaconda3\lib\site-packages (from scikit-optimize)

Collecting pyaml>=16.9 (from scikit-optimize)

Downloading pyaml-25.1.0-py3-none-any.whl.metadata (12 kB)

Requirement already satisfied: numpy>=1.20.3 in c:\users\lsi8012\appdata\local\anaconda3\lib\site-packages (from pyaml>=16.9)

Requirement already satisfied: scipy>=1.1.0 in c:\users\lsi8012\appdata\local\anaconda3\lib\site-packages (from pyaml>=16.9)

Requirement already satisfied: scikit-learn>=1.0.0 in c:\users\lsi8012\appdata\local\anaconda3\lib\site-packages (from pyaml>=16.9)

Requirement already satisfied: packaging>=21.3 in c:\users\lsi8012\appdata\roaming\python\python39\lib\site-packages (from pyaml>=16.9)

Requirement already satisfied: PyYAML in c:\users\lsi8012\appdata\local\anaconda3\lib\site-packages (from pyaml>=16.9)

Requirement already satisfied: threadpoolctl>=3.1.0 in c:\users\lsi8012\appdata\local\anaconda3\lib\site-packages (from pyaml>=16.9)

Downloading scikit\_optimize-0.10.2-py2.py3-none-any.whl (107 kB)

----- 0.0/107.8 kB ? eta -:--:--

----- -- 102.4/107.8 kB 5.8 MB/s eta 0:00:01

----- 107.8/107.8 kB 3.1 MB/s eta 0:00:00

Downloading pyaml-25.1.0-py3-none-any.whl (26 kB)

Installing collected packages: pyaml, scikit-optimize

Successfully installed pyaml-25.1.0 scikit-optimize-0.10.2

```
from skopt import BayesSearchCV
from skopt.space import Integer, Categorical
```

```
# Step 3: Define search space for Bayesian optimization
```

```
search_space = {
```

```
    'knn__n_neighbors': Integer(1, 19), # Odd values will be sampled if needed
```

```
    'knn__weights': Categorical(['uniform', 'distance']),
```

```
    'knn__metric': Categorical(['euclidean', 'manhattan', 'minkowski']),
```

```
    'knn__p': Integer(1, 2) # Used only when metric is minkowski
```

```
}
```

```
# Step 4: Set up BayesSearchCV
bayes_search = BayesSearchCV(
    estimator=pipeline,
    search_spaces=search_space,
    n_iter=30, # Number of different combinations to try
    scoring='neg_root_mean_squared_error',
    cv=5,
    n_jobs=-1,
    verbose=1,
    random_state=42
)
```

```
# Step 5: Fit BayesSearchCV
print("Tuning KNN hyperparameters with Bayesian Optimization...")
bayes_search.fit(X_train, y_train)

# Get best parameters and best score
best_params = bayes_search.best_params_
best_score = -bayes_search.best_score_ # Convert negative RMSE to positive

# Display results
print("\nBest Parameters (Bayesian Optimization):")
for param, value in best_params.items():
    print(f"{param}: {value}")
print(f"Best CV RMSE: {best_score:.2f}")
```

[illegible]

Fitting 5 folds for each of 1 candidates, totalling 5 fits  
Fitting 5 folds for each of 1 candidates, totalling 5 fits  
Fitting 5 folds for each of 1 candidates, totalling 5 fits

```
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\skopt\optimizer\optimizer.py:517:  
    warnings.warn(  

```

Fitting 5 folds for each of 1 candidates, totalling 5 fits

```
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\skopt\optimizer\optimizer.py:517:  
    warnings.warn(  

```

Fitting 5 folds for each of 1 candidates, totalling 5 fits

```
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\skopt\optimizer\optimizer.py:517:  
    warnings.warn(  

```

Fitting 5 folds for each of 1 candidates, totalling 5 fits  
Fitting 5 folds for each of 1 candidates, totalling 5 fits  
Fitting 5 folds for each of 1 candidates, totalling 5 fits  
Fitting 5 folds for each of 1 candidates, totalling 5 fits

```
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\skopt\optimizer\optimizer.py:517:  
    warnings.warn(  

```

Fitting 5 folds for each of 1 candidates, totalling 5 fits

```
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\skopt\optimizer\optimizer.py:517:  
    warnings.warn(  

```

Fitting 5 folds for each of 1 candidates, totalling 5 fits

```
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\skopt\optimizer\optimizer.py:517:  
    warnings.warn(  

```

Fitting 5 folds for each of 1 candidates, totalling 5 fits

```
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\skopt\optimizer\optimizer.py:517:  
    warnings.warn(  

```



Fitting 5 folds for each of 1 candidates, totalling 5 fits

```
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\skopt\optimizer\optimizer.py:517:
  warnings.warn(
```

Fitting 5 folds for each of 1 candidates, totalling 5 fits

Fitting 5 folds for each of 1 candidates, totalling 5 fits

Best Parameters (Bayesian Optimization):

knn\_\_metric: manhattan

knn\_\_n\_neighbors: 8

knn\_\_p: 2

knn\_\_weights: distance

Best CV RMSE: 4005.70

```
# Step 6: Evaluate on test set
best_model = bayes_search.best_estimator_
y_pred = best_model.predict(X_test)

# Calculate RMSE on test set
test_rmse = root_mean_squared_error(y_test, y_pred)
print(f"Test RMSE: {test_rmse:.2f}")
```

Test RMSE: 3811.89

## 3 Hyperparameter tuning

In this chapter we'll introduce several functions that help with tuning hyperparameters of a machine learning model.

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split, cross_val_score, cross_val_predict, \
cross_validate, GridSearchCV, RandomizedSearchCV, KFold, StratifiedKFold, RepeatedKFold, Rep
from sklearn.neighbors import KNeighborsClassifier, KNeighborsRegressor
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, recall_score, mean_squared_error
from scipy.stats import uniform
from skopt import BayesSearchCV
from skopt.space import Real, Categorical, Integer
import seaborn as sns
from skopt.plots import plot_objective, plot_histogram, plot_convergence
import matplotlib.pyplot as plt
import warnings
from IPython import display
```

Let us read and pre-process data first. Then we'll be ready to tune the model hyperparameters. We'll use KNN as the model. Note that KNN has multiple hyperparameters to tune, such as number of neighbors, distance metric, weights of neighbours, etc.

```
#Using the same datasets as used for linear regression in STAT303-2,
#so that we can compare the non-linear models with linear regression
trainf = pd.read_csv('./Datasets/Car_features_train.csv')
trainp = pd.read_csv('./Datasets/Car_prices_train.csv')
testf = pd.read_csv('./Datasets/Car_features_test.csv')
testp = pd.read_csv('./Datasets/Car_prices_test.csv')
train = pd.merge(trainf, trainp)
test = pd.merge(testf, testp)
train.head()
```

	carID	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	18473	bmw	6 Series	2020	Semi-Auto	11	Diesel	145	53.3282	3.0	37980
1	15064	bmw	6 Series	2019	Semi-Auto	10813	Diesel	145	53.0430	3.0	33980
2	18268	bmw	6 Series	2020	Semi-Auto	6	Diesel	145	53.4379	3.0	36850
3	18480	bmw	6 Series	2017	Semi-Auto	18895	Diesel	145	51.5140	3.0	25998
4	18492	bmw	6 Series	2015	Automatic	62953	Diesel	160	51.4903	3.0	18990

```

predictors = ['mpg', 'engineSize', 'year', 'mileage']
X_train = train[predictors]
y_train = train['price']
X_test = test[predictors]
y_test = test['price']

# Scale
sc = StandardScaler()

sc.fit(X_train)
X_train_scaled = sc.transform(X_train)
X_test_scaled = sc.transform(X_test)

```

### 3.1 GridSearchCV

The function is used to compute the cross-validated score (*MSE*, *RMSE*, *accuracy*, *etc.*) over a grid of hyperparameter values. This helps avoid nested `for()` loops if multiple hyperparameter values need to be tuned.

```

# GridSearchCV works in three steps:

# 1) Create the model
model = KNeighborsRegressor() # No inputs defined inside the model

# 2) Create a hyperparameter grid (as a dict)
# the keys should be EXACTLY the same as the names of the model inputs
# the values should be an array or list of hyperparam values you want to try out

# 30 K values x 2 weight settings x 3 metric settings = 180 different combinations in this grid
grid = {'n_neighbors': np.arange(5, 151, 5), 'weights':['uniform', 'distance'],
        'metric': ['manhattan', 'euclidean', 'chebyshev']}

# 3) Create the Kfold object (Using RepeatedKfold will be more robust, but more expensive, use it if you
# have the budget)

```

```

kfold = KFold(n_splits = 5, shuffle = True, random_state = 1)

# 4) Create the CV object
# Look at the documentation to see the order in which the objects must be specified within the
gcv = GridSearchCV(model, grid, cv = kfold, scoring = 'neg_root_mean_squared_error', n_jobs=5)

# Fit the models, and cross-validate
gcv.fit(X_train_scaled, y_train)

```

Fitting 5 folds for each of 180 candidates, totalling 900 fits

```

GridSearchCV(cv=KFold(n_splits=5, random_state=1, shuffle=True),
             estimator=KNeighborsRegressor(), n_jobs=-1,
             param_grid={'metric': ['manhattan', 'euclidean', 'chebyshev'],
                         'n_neighbors': array([ 5, 10, 15, 20, 25, 30, 35, 40, 45,
70, 75, 80, 85, 90, 95, 100, 105, 110, 115, 120, 125, 130,
135, 140, 145, 150]),
                         'weights': ['uniform', 'distance']}},
             scoring='neg_root_mean_squared_error', verbose=10)

```

The optimal estimator based on cross-validation is:

```
gcv.best_estimator_
```

```
KNeighborsRegressor(metric='manhattan', n_neighbors=10, weights='distance')
```

The optimal hyperparameter values (*based on those considered in the grid search*) are:

```
gcv.best_params_
```

```
{'metric': 'manhattan', 'n_neighbors': 10, 'weights': 'distance'}
```

The cross-validated root mean squared error for the optimal hyperparameter values is:

```
-gcv.best_score_
```

```
5740.928686723918
```

The RMSE on test data for the optimal hyperparameter values is:

```
y_pred = gcv.predict(X_test_scaled)
mean_squared_error(y_test, y_pred, squared=False)
```

5747.466851437544

Note that the error is further reduced as compared to the case when we tuned only one hyperparameter in the [previous chapter](#). We must tune all the hyperparameters that can effect prediction accuracy, in order to get the most accurate model.

The results for each cross-validation are stored in the `cv_results_` attribute.

```
pd.DataFrame(gcv.cv_results_).head()
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_metric	param_n_neighb
0	0.011169	0.005060	0.011768	0.001716	manhattan	5
1	0.009175	0.001934	0.009973	0.000631	manhattan	5
2	0.008976	0.001092	0.012168	0.001323	manhattan	10
3	0.007979	0.000001	0.011970	0.000892	manhattan	10
4	0.006781	0.000748	0.012367	0.001017	manhattan	15

These results can be useful to see if other hyperparameter values are almost equally good.

For example, the next two best optimal values of the hyperparameter correspond to neighbors being 15 and 5 respectively. As the test error has a high variance, the best hyperparameter values need not necessarily be actually optimal.

```
pd.DataFrame(gcv.cv_results_).sort_values(by = 'rank_test_score').head()
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_metric	param_n_neighb
3	0.007979	0.000001	0.011970	0.000892	manhattan	10
5	0.009374	0.004829	0.013564	0.001850	manhattan	15
1	0.009175	0.001934	0.009973	0.000631	manhattan	5
7	0.007977	0.001092	0.017553	0.002054	manhattan	20
9	0.007777	0.000748	0.019349	0.003374	manhattan	25

Let us compute the RMSE on test data based on the 2nd and 3rd best hyperparameter values.

```
model = KNeighborsRegressor(n_neighbors=15, metric='manhattan', weights='distance').fit(X_train_scaled, y_train)
mean_squared_error(model.predict(X_test_scaled), y_test, squared = False)
```

5800.418957612656

```
model = KNeighborsRegressor(n_neighbors=5, metric='manhattan', weights='distance').fit(X_train_scaled, y_train)
mean_squared_error(model.predict(X_test_scaled), y_test, squared = False)
```

5722.4859230146685

We can see that the RMSE corresponding to the 3rd best hyperparameter value is the least. Due to variance in test errors, it may be a good idea to consider the set of top few best hyperparameter values, instead of just considering the best one.

## 3.2 RandomizedSearchCV()

In case of many possible values of hyperparameters, it may be computationally very expensive to use GridSearchCV(). In such cases, RandomizedSearchCV() can be used to compute the cross-validated score on a randomly selected subset of hyperparameter values from the specified grid. The number of values can be fixed by the user, as per the available budget.

```
# RandomizedSearchCV works in three steps:

# 1) Create the model
model = KNeighborsRegressor() # No inputs defined inside the model

# 2) Create a hyperparameter grid (as a dict)
# the keys should be EXACTLY the same as the names of the model inputs
# the values should be an array or list of hyperparam values, or distribution of hyperparameters

grid = {'n_neighbors': range(1, 500), 'weights': ['uniform', 'distance'],
        'metric': ['minkowski'], 'p': uniform(loc=1, scale=10)} #We can specify a distribution
                                                                #for continuous hyperparameters

# 3) Create the Kfold object (Using RepeatedKfold will be more robust, but more expensive, use it if budget allows)
# have the budget)
kfold = KFold(n_splits = 5, shuffle = True, random_state = 1)
```

```
# 4) Create the CV object
# Look at the documentation to see the order in which the objects must be specified within the
gcv = RandomizedSearchCV(model, param_distributions = grid, cv = kfold, n_iter = 180, random_state=10,
                        scoring = 'neg_root_mean_squared_error', n_jobs = -1, verbose = 10)

# Fit the models, and cross-validate
gcv.fit(X_train_scaled, y_train)
```

Fitting 5 folds for each of 180 candidates, totalling 900 fits

```
RandomizedSearchCV(cv=KFold(n_splits=5, random_state=1, shuffle=True),
                  estimator=KNeighborsRegressor(), n_iter=180, n_jobs=-1,
                  param_distributions={'metric': ['minkowski'],
                                       'n_neighbors': range(1, 500),
                                       'p': <scipy.stats._distn_infrastructure.rv_continuous object>,
                                       'weights': ['uniform', 'distance']},
                  random_state=10, scoring='neg_root_mean_squared_error',
                  verbose=10)
```

```
gcv.best_params_
```

```
{'metric': 'minkowski',
 'n_neighbors': 3,
 'p': 1.252639454318171,
 'weights': 'uniform'}
```

```
gcv.best_score_
```

```
-6239.171627183809
```

```
y_pred = gcv.predict(X_test_scaled)
mean_squared_error(y_test, y_pred, squared=False)
```

```
6176.533397589911
```

Note that in this example, `RandomizedSearchCV()` helps search for optimal values of the hyperparameter  $p$  over a continuous domain space. In this dataset,  $p = 1$  seems to be the optimal value. However, if the optimal value was somewhere in the middle of a larger

continuous domain space (*instead of the boundary of the domain space*), and there were several other hyperparameters, some of which were not influencing the response (*effect sparsity*), `RandomizedSearchCV()` is likely to be more effective in estimating the optimal value of the continuous hyperparameter.

The advantages of `RandomizedSearchCV()` over `GridSearchCV()` are:

1. `RandomizedSearchCV()` fixes the computational cost in case of large number of hyperparameters / large number of levels of individual hyperparameters. If there are  $n$  hyperparameters, each with 3 levels, the number of all possible hyperparameter values will be  $3^n$ . The computational cost increase exponentially with increase in number of hyperparameters.
2. In case of a hyperparameter having continuous values, the distribution of the hyperparameter can be specified in `RandomizedSearchCV()`.
3. In case of effect sparsity of hyperparameters, i.e., if only a few hyperparameters significantly effect prediction accuracy, `RandomizedSearchCV()` is likely to consider more unique values of the influential hyperparameters as compared to `GridSearchCV()`, and is thus likely to provide more optimal hyperparameter values as compared to `GridSearchCV()`. The figure below shows effect sparsity where there are 2 hyperparameters, but only one of them is associated with the cross-validated score, Here, it is more likely that the optimal cross-validated score will be obtained by `RandomizedSearchCV()`, as it is evaluating the model on 9 unique values of the relevant hyperparameter, instead of just 3.

<IPython.core.display.Image object>

### 3.3 `BayesSearchCV()`

Unlike the grid search and random search, which treat hyperparameter sets independently, the Bayesian optimization is an informed search method, meaning that it learns from previous iterations. The number of trials in this approach is determined by the user.

- The function begins by computing the cross-validated score by randomly selecting a few hyperparameter values from the specified distribution of hyperparameter values.
- Based on the data of hyperparameter values tested (*predictors*), and the cross-validated score (*the response*), a Gaussian process model is developed to estimate the cross-validated score & the uncertainty in the estimate in the entire space of the hyperparameter values



- A criterion that “explores” uncertain regions of the space of hyperparameter values (*where it is difficult to predict cross-validated score*), and “exploits” promising regions of the space of hyperparameter values (*where the cross-validated score is predicted to minimize*) is used to suggest the next hyperparameter value that will potentially minimize the cross-validated score
- Cross-validated score is computed at the suggested hyperparameter value, the Gaussian process model is updated, and the previous step is repeated, until a certain number of iterations specified by the user.

To summarize, instead of blindly testing the model for the specified hyperparameter values (*as in `GridSearchCV()`*), or randomly testing the model on certain hyperparameter values (*as in `RandomizedSearchCV()`*), `BayesSearchCV()` smartly tests the model for those hyperparameter values that are likely to reduce the cross-validated score. The algorithm becomes “smarter” as it “learns” more with increasing iterations.

Here is a nice [blog](#), if you wish to understand more about the Bayesian optimization procedure.

```
# BayesSearchCV works in three steps:

# 1) Create the model
model = KNeighborsRegressor(metric = 'minkowski') # No inputs defined inside the model

# 2) Create a hyperparameter grid (as a dict)
# the keys should be EXACTLY the same as the names of the model inputs
# the values should be the distribution of hyperparameter values. Lists and NumPy arrays can
# also be used

grid = {'n_neighbors': Integer(1, 500), 'weights': Categorical(['uniform', 'distance']),
        'p': Real(1, 10, prior = 'uniform')}

# 3) Create the Kfold object (Using RepeatedKfold will be more robust, but more expensive,
# use it if you have the budget)
kfold = KFold(n_splits = 5, shuffle = True, random_state = 1)

# 4) Create the CV object
# Look at the documentation to see the order in which the objects must be specified within
# the function
gcv = BayesSearchCV(model, search_spaces = grid, cv = kfold, n_iter = 180, random_state = 10,
                    scoring = 'neg_root_mean_squared_error', n_jobs = -1)

# Fit the models, and cross-validate
```

```
# Sometimes the Gaussian process model predicting the cross-validated score suggests a
# "promising point" (i.e., set of hyperparameter values) for cross-validation that it has
# already suggested earlier. In such a case a warning is raised, and the objective
# function (i.e., the cross-validation score) is computed at a randomly selected point
# (as in RandomizedSearchCV()). This feature helps the algorithm explore other regions of
# the hyperparameter space, rather than only searching in the promising regions. Thus, it
# balances exploration (of the hyperparameter space) with exploitation (of the promising
# regions of the hyperparameter space)

warnings.filterwarnings("ignore")
gcv.fit(X_train_scaled, y_train)
warnings.resetwarnings()
```

The optimal hyperparameter values (*based on Bayesian search*) on the provided distribution of hyperparameter values are:

```
gcv.best_params_
```

```
OrderedDict([('n_neighbors', 9),
             ('p', 1.0008321732366932),
             ('weights', 'distance')])
```

The cross-validated root mean squared error for the optimal hyperparameter values is:

```
-gcv.best_score_
```

```
5756.172382596493
```

The RMSE on test data for the optimal hyperparameter values is:

```
y_pred = gcv.predict(X_test_scaled)
mean_squared_error(y_test, y_pred, squared=False)
```

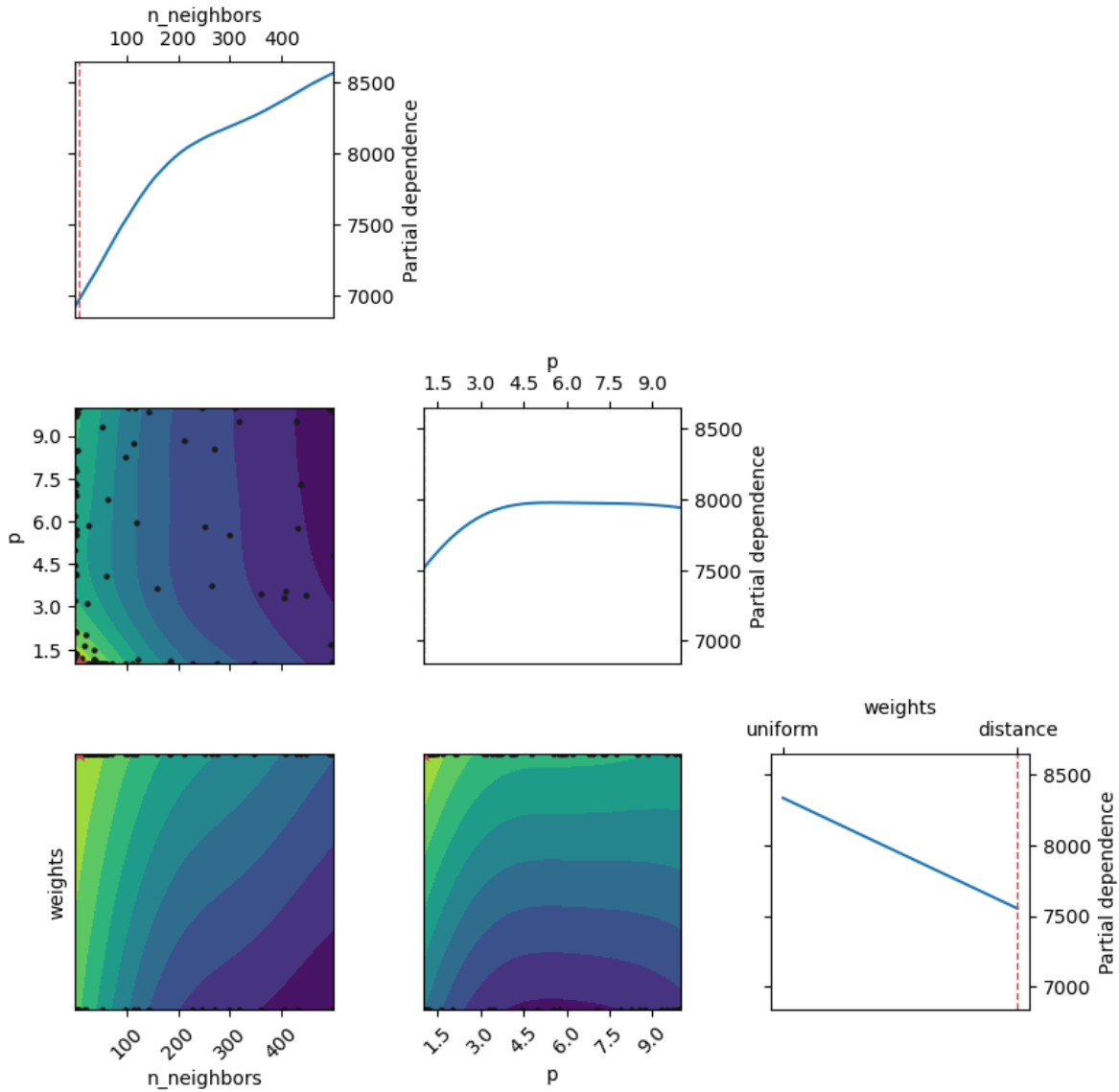
```
5740.432278861367
```

### 3.3.1 Diagnosis of cross-validated score optimization

Below are the partial dependence plots of the objective function (*i.e.*, the cross-validated score). The cross-validated score predictions are based on the most recently updated model (*i.e.*, the updated *Gaussian Process* model at the end of `n_iter` iterations specified by the user) that predicts the cross-validated score.

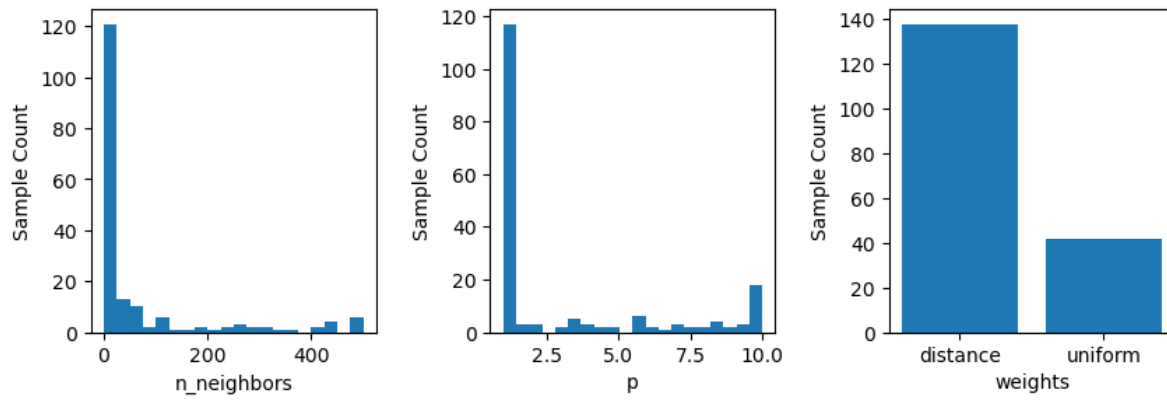
Check the `plot_objective()` documentation to interpret the plots.

```
plot_objective(gcv.optimizer_results_[0],  
               dimensions=["n_neighbors", "p", "weights"], size = 3)  
plt.show();
```



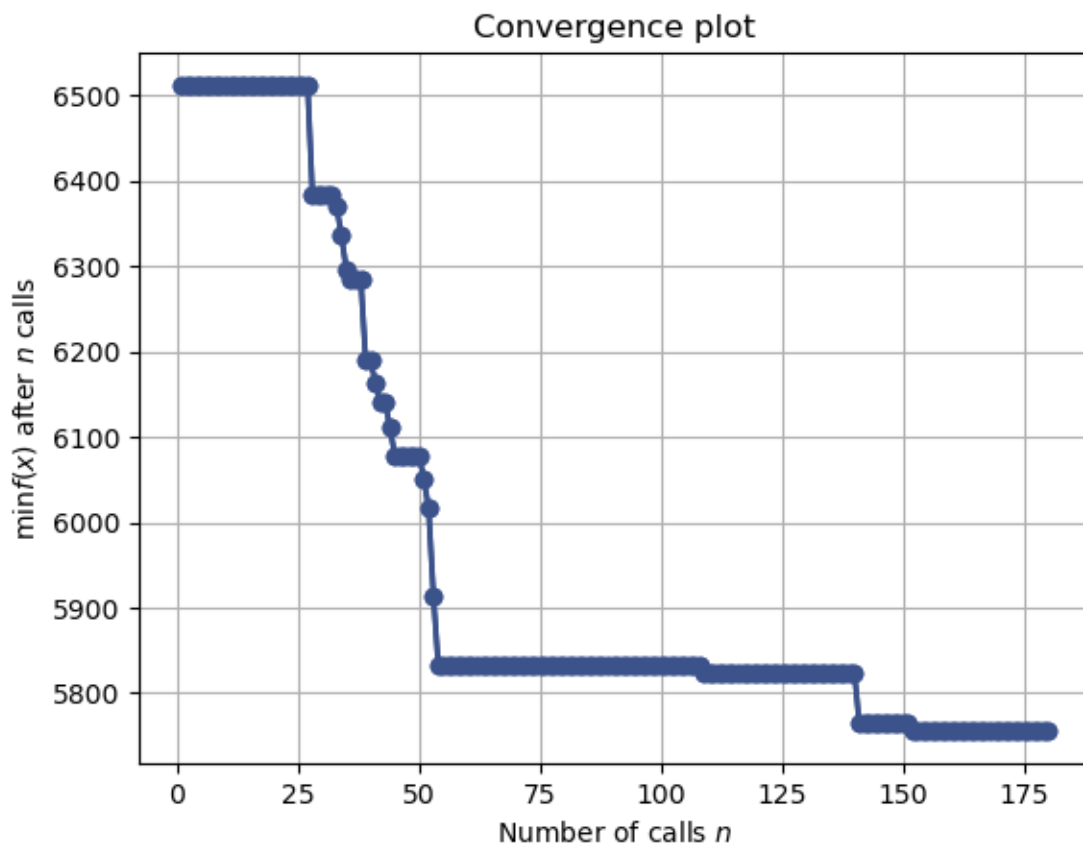
The frequency of individual hyperparameter values considered can also be visualized as below.

```
fig, ax = plt.subplots(1, 3, figsize = (10, 3))
plt.subplots_adjust(wspace=0.4)
plot_histogram(gcv.optimizer_results_[0], 0, ax = ax[0])
plot_histogram(gcv.optimizer_results_[0], 1, ax = ax[1])
plot_histogram(gcv.optimizer_results_[0], 2, ax = ax[2])
plt.show()
```



Below is the plot showing the minimum cross-validated score computed obtained until 'n' hyperparameter values are considered for cross-validation.

```
plot_convergence(gcv.optimizer_results_)
plt.show()
```



Note that the cross-validated error is close to the optimal value in the 53rd iteration itself.

The cross-validated error at the 53rd iteration is:

```
gcv.optimizer_results_[0]['func_vals'][53]
```

```
5831.87280274334
```

The hyperparameter values at the 53rd iterations are:

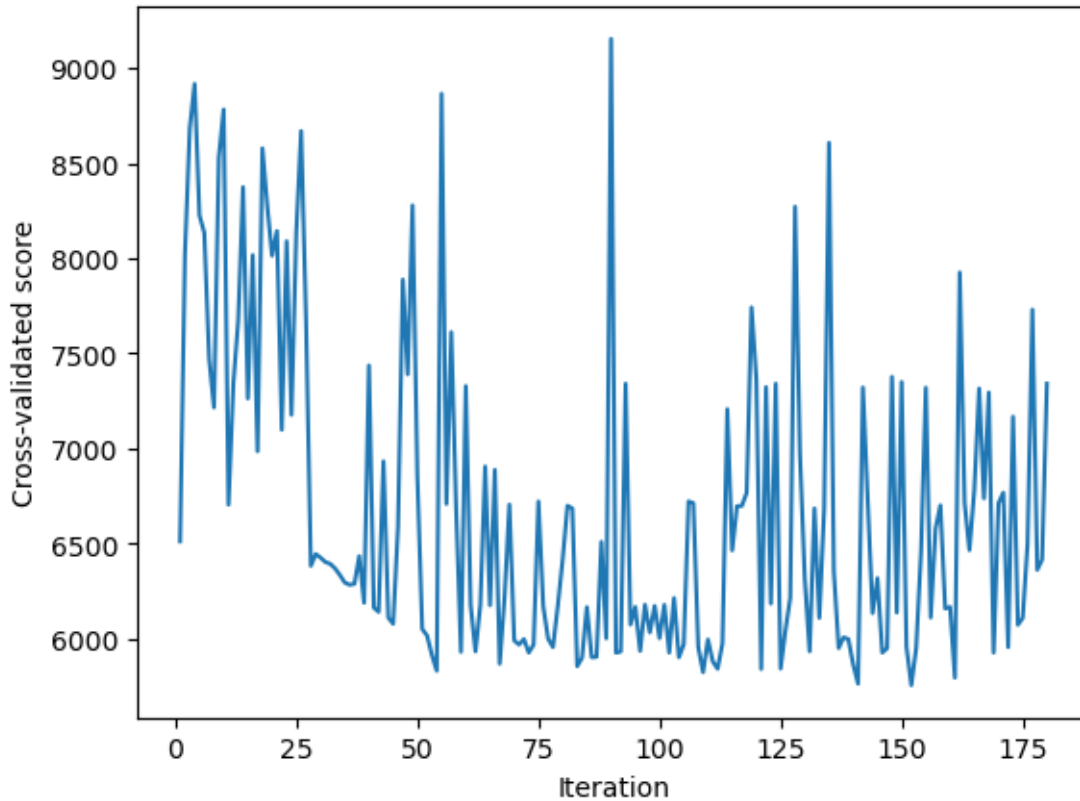
```
gcv.optimizer_results_[0]['x_iters'][53]
```

```
[15, 1.0, 'distance']
```

Note that this is the 2nd most optimal hyperparameter value based on `GridSearchCV()`.

Below is the plot showing the cross-validated score computed at each of the 180 hyperparameter values considered for cross-validation. The plot shows that the algorithm seems to explore new regions of the domain space, instead of just exploiting the promising ones. There is a balance between exploration and exploitation for finding the optimal hyperparameter values that minimize the objective function (*i.e., the function that models the cross-validated score*).

```
sns.lineplot(x = range(1, 181), y = gcv.optimizer_results_[0]['func_vals'])
plt.xlabel('Iteration')
plt.ylabel('Cross-validated score')
plt.show();
```



The advantages of `BayesSearchCV()` over `GridSearchCV()` and `RandomizedSearchCV()` are:

1. The Bayesian Optimization approach gives the benefit that we can give a much larger range of possible values, since over time we identify and exploit the most promising regions and discard the not so promising ones. Plain grid-search would burn computational resources to explore all regions of the domain space with the same granularity, even the not promising ones. Since we search much more effectively in Bayesian search, we can search over a larger domain space.
2. BayesSearch CV may help us identify the optimal hyperparameter value in fewer iterations if the Gaussian process model estimating the cross-validated score is relatively accurate. However, this is not certain. Grid and random search are completely uninformed by past evaluations, and as a result, often spend a significant amount of time evaluating “bad” hyperparameters.
3. BayesSearch CV is more reliable in cases of a large search space, where random selection may miss sampling values from optimal regions of the search space.

The disadvantages of `BayesSearchCV()` over `GridSearchCV()` and `RandomizedSearchCV()` are:

1. `BayesSearchCV()` has a cost of learning from past data, i.e., updating the model that predicts the cross-validated score after every iteration of evaluating the cross-validated score on a new hyperparameter value. This cost will continue to increase as more and more data is collected. There is no such cost in `GridSearchCV()` and `RandomizedSearchCV()` as there is no learning. This implies that each iteration of `BayesSearchCV()` will take a longer time than each iteration of `GridSearchCV()` / `RandomizedSearchCV()`. Thus, even if `BayesSearchCV()` finds the optimal hyperparameter value in fewer iterations, it may take more time than `GridSearchCV()` / `RandomizedSearchCV()` for the same.
2. The success of `BayesSearchCV()` depends on the predictions and associated uncertainty estimated by the Gaussian process (GP) model that predicts the cross-validated score. The GP model, although works well in general, may not be suitable for certain datasets, or may take a relatively large number of iterations to learn for certain datasets.

### 3.3.2 Live monitoring of cross-validated score

Note that it will be useful monitor the cross-validated score while the Bayesian Search CV code is running, and stop the code as soon as the desired accuracy is reached, or the optimal cross-validated score doesn't seem to improve. The `fit()` method of the `BayesSearchCV()` object has a `callback` argument that can be used as follows:

```
model = KNeighborsRegressor(metric = 'minkowski') # No inputs defined inside the model
grid = {'n_neighbors': Integer(1, 500), 'weights': Categorical(['uniform', 'distance']),
        'p': Real(1, 10, prior = 'uniform')}

kfold = KFold(n_splits = 5, shuffle = True, random_state = 1)
gcv = BayesSearchCV(model, search_spaces = grid, cv = kfold, n_iter = 180, random_state = 10,
                    scoring = 'neg_root_mean_squared_error', n_jobs = -1)

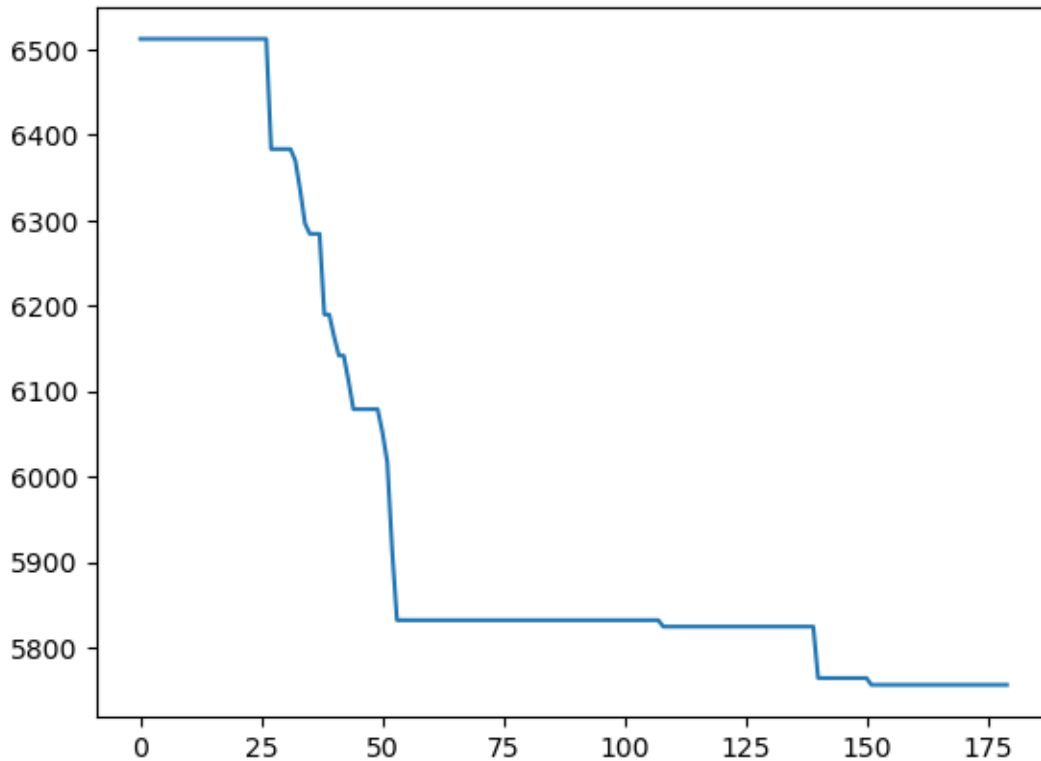
paras = list(gcv.search_spaces.keys())
paras.sort()

def monitor(optim_result):
    cv_values = pd.Series(optim_result['func_vals']).cummin()
    display.clear_output(wait = True)
    min_ind = pd.Series(optim_result['func_vals']).argmin()
    print(paras, "=", optim_result['x_iters'][min_ind], pd.Series(optim_result['func_vals']))
    sns.lineplot(cv_values)
    plt.show()
```



```
gcv.fit(X_train_scaled, y_train, callback = monitor)
```

```
['n_neighbors', 'p', 'weights'] = [9, 1.0008321732366932, 'distance'] 5756.172382596493
```



```
BayesSearchCV(cv=KFold(n_splits=5, random_state=1, shuffle=True),
               estimator=KNeighborsRegressor(), n_iter=180, n_jobs=-1,
               random_state=10, scoring='neg_root_mean_squared_error',
               search_spaces={'n_neighbors': Integer(low=1, high=500, prior='uniform', transform='log'),
                              'p': Real(low=1, high=10, prior='uniform', transform='normalize'),
                              'weights': Categorical(categories=('uniform', 'distance'), prior='uniform')})
```

### 3.4 cross\_validate()

We have used `cross_val_score()` and `cross_val_predict()` so far.

When can we use one over the other?

The function `cross_validate()` is similar to `cross_val_score()` except that it has the option to return multiple cross-validated metrics, instead of a single one.

Consider the heart disease classification problem, where the response is **target** (*whether the person has a heart disease or not*).

```
data = pd.read_csv('Datasets/heart_disease_classification.csv')
data.head()
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

Let us pre-process the data.

```
# First, separate the response and the predictors
y = data['target']
X = data.drop('target', axis=1)
```

```
# Separate the data (X,y) into training and test
```

```
# Inputs:
# data
# train-test ratio
# random_state for reproducible code
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=20, stratify=y)
```

```
# stratify=y makes sure the class 0 to class 1 ratio in the training and test sets are kept the same
```

```
model = KNeighborsClassifier()
sc = StandardScaler()
sc.fit(X_train)
X_train_scaled = sc.transform(X_train)
X_test_scaled = sc.transform(X_test)
```

Suppose we want to take recall above a certain threshold with the highest precision possible. `cross_validate()` computes the cross-validated score for multiple metrics - rest is the same as `cross_val_score()`.

```

Ks = np.arange(10,200,10)

scores = []

for K in Ks:
    model = KNeighborsClassifier(n_neighbors=K) # Keeping distance uniform
    scores.append(cross_validate(model, X_train_scaled, y_train, cv=5, scoring = ['accuracy'

scores

# The output is now a list of dicts - easy to convert to a df

df_scores = pd.DataFrame(scores) # We need to handle test_recall and test_precision cols

df_scores['CV_recall'] = df_scores['test_recall'].apply(np.mean)
df_scores['CV_precision'] = df_scores['test_precision'].apply(np.mean)
df_scores['CV_accuracy'] = df_scores['test_accuracy'].apply(np.mean)

df_scores.index = Ks # We can set K values as indices for convenience

#df_scores
# What happens as K increases?
# Recall increases (not monotonically)
# Precision decreases (not monotonically)
# Why?
# Check the class distribution in the data - more obs with class 1
# As K gets higher, the majority class overrules (visualized in the slides)
# More 1s means less FNs - higher recall
# More 1s means more FPs - lower precision
# Would this be the case for any dataset?
# NO!! Depends on what the majority class is!

```

Suppose we wish to have the maximum possible precision for at least 95% recall.

The optimal  $K$  will be:

```
df_scores.loc[df_scores['CV_recall'] > 0.95, 'CV_precision'].idxmax()
```

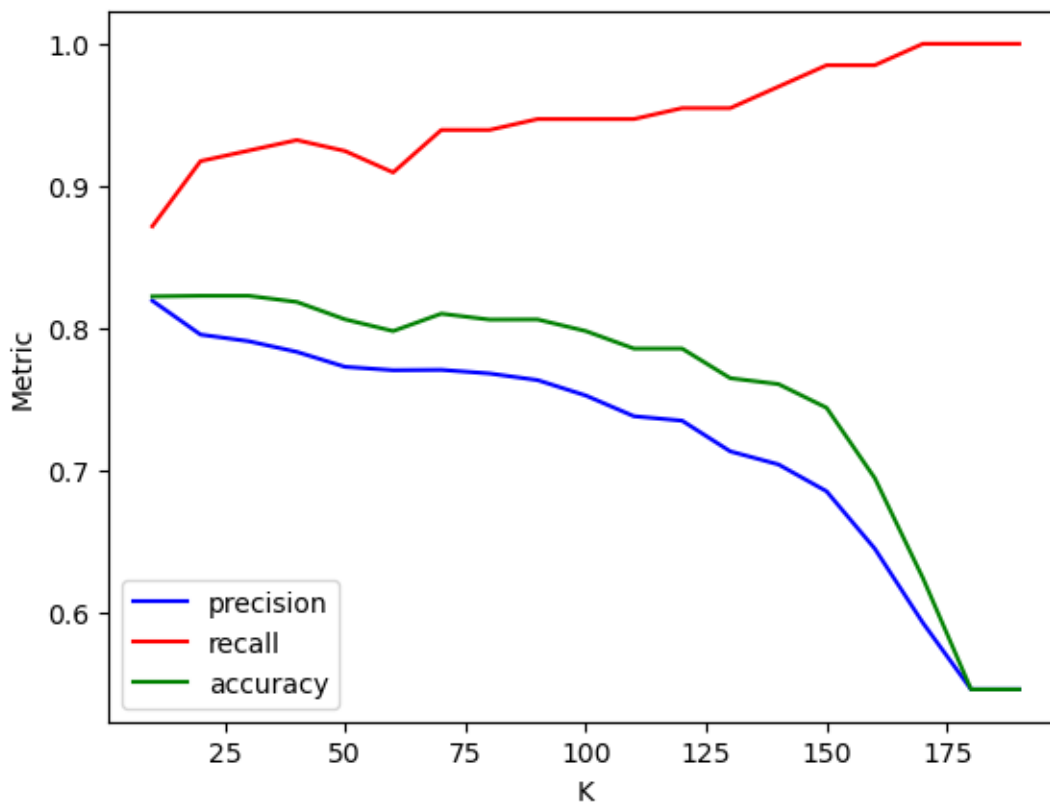
120

The cross-validated precision, recall and accuracy for the optimal  $K$  are:

```
df_scores.loc[120, ['CV_recall', 'CV_precision', 'CV_accuracy']]
```

```
CV_recall      0.954701  
CV_precision    0.734607  
CV_accuracy    0.785374  
Name: 120, dtype: object
```

```
sns.lineplot(x = df_scores.index, y = df_scores.CV_precision, color = 'blue', label = 'precision')  
sns.lineplot(x = df_scores.index, y = df_scores.CV_recall, color = 'red', label = 'recall')  
sns.lineplot(x = df_scores.index, y = df_scores.CV_accuracy, color = 'green', label = 'accuracy')  
plt.ylabel('Metric')  
plt.xlabel('K')  
plt.show()
```



## **Part II**

# **Tree based models**

## 4 Regression trees

*Read section 8.1.1 of the book before using these notes.*

*Note that in this course, lecture notes are not sufficient, you must read the book for better understanding. Lecture notes are just implementing the concepts of the book on a dataset, but not explaining the concepts elaborately.*

```
# Importing necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(font_scale=1.35)

# import the decision tree regressor
from sklearn.tree import DecisionTreeRegressor, plot_tree, export_graphviz

# split the dataset into training and testing sets
from sklearn.model_selection import train_test_split

from sklearn.model_selection import cross_val_score, GridSearchCV, cross_val_predict, KFold

from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, FunctionTransformer
from sklearn.metrics import root_mean_squared_error, r2_score, make_scorer
```

We will use the same dataset as in the KNN model for regression trees.

```
# Load the dataset
car = pd.read_csv('Datasets/car.csv')
car.head()
```

	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	vw	Beetle	2014	Manual	55457	Diesel	30	65.3266	1.6	7490
1	vauxhall	GTC	2017	Manual	15630	Petrol	145	47.2049	1.4	10998
2	merc	G Class	2012	Automatic	43000	Diesel	570	25.1172	3.0	44990
3	audi	RS5	2019	Automatic	10	Petrol	145	30.5593	2.9	51990
4	merc	X-CLASS	2018	Automatic	14000	Diesel	240	35.7168	2.3	28990

## 4.1 Native Support for Missing Values

Starting with **scikit-learn version 1.3**, classical tree-based models in **scikit-learn** have added **native support for missing values**, which simplifies preprocessing and improves model robustness:

- `DecisionTreeClassifier` supports missing values as of **version 1.3.0**
- `RandomForestClassifier` adds support in **version 1.4.0**

This means you no longer need to impute missing values manually before training these models.

To take advantage of this feature, first check your **scikit-learn** version:

```
import sklearn
print(sklearn.__version__)
```

1.6.1

If your version is below 1.4.0, you can upgrade by running:

```
# pip install --upgrade scikit-learn
```

```
# Make a copy of the original dataset
car_missing = car.copy()

# Randomly add missing values
# Inject missing values into 10% of the 'mileage' column
car_missing.loc[car_missing.sample(frac=0.1, random_state=42).index, 'mileage'] = np.nan
# Inject missing values into 10% of the 'fuelType' and 'engineSize' columns

car_missing.loc[car_missing.sample(frac=0.1, random_state=42).index, 'fuelType'] = np.nan
car_missing.loc[car_missing.sample(frac=0.1, random_state=42).index, 'engineSize'] = np.nan
```

```
car_missing.isna().sum()
```

```
brand          0
model          0
year           0
transmission   0
mileage        763
fuelType       763
tax            0
mpg            0
engineSize     763
price          0
dtype: int64
```

```
# Split the car_missing dataset into features and target
X_missing = car_missing.drop(columns=['price'])
y_missing = car_missing['price']
```

#### 4.1.1 Build a regression tree using mileage as the solo predictor

```
# Use only 'mileage' as the feature
X_mileage = X_missing[['mileage']]
y_mileage = y_missing
```

```
# Create a DecisionTreeRegressor model
reg_tree = DecisionTreeRegressor(random_state=42)

# Fit the model to the data
reg_tree.fit(X_mileage, y_mileage)
```

```
DecisionTreeRegressor(random_state=42)
```

```
# Predict the target variable using the model
y_pred = reg_tree.predict(X_mileage)

# Calculate the RMSE and R2 score
rmse = np.sqrt(np.mean((y_missing - y_pred) ** 2))
r2 = r2_score(y_missing, y_pred)
print(f"RMSE: {rmse:.2f}")
print(f"R2: {r2:.2f}")
```



RMSE: 8797.85

$R^2$ : 0.71

## 4.2 Building regression trees

```
X = car.drop(columns=['price'])
y = car['price']
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

### 4.2.1 Using only mileage feature

```
# Use only 'mileage' as the feature
X_train_mileage = X_train[['mileage']]
X_test_mileage = X_test[['mileage']]

# Create a DecisionTreeRegressor model
reg_tree = DecisionTreeRegressor(random_state=42, max_depth=3)

# Fit the model to the training data
reg_tree.fit(X_train_mileage, y_train)

# Predict the target variable using the model
y_pred = reg_tree.predict(X_test_mileage)

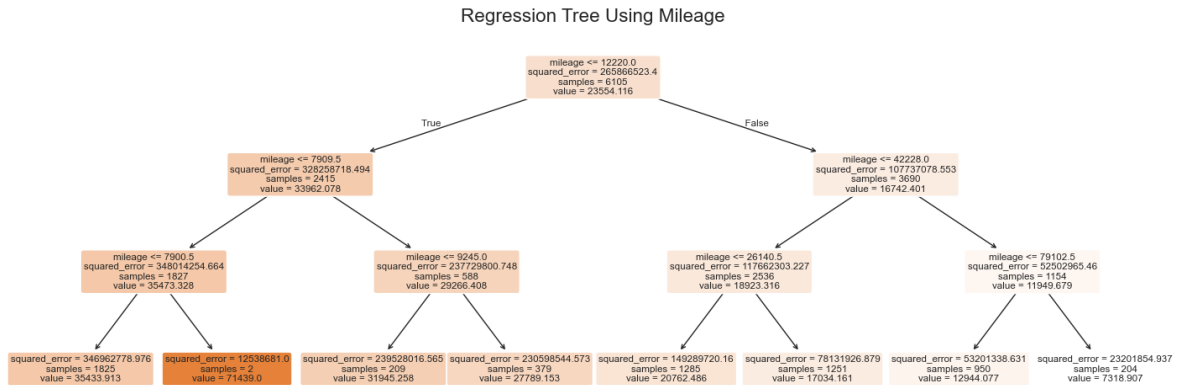
# Calculate the RMSE and  $R^2$  score
rmse = np.sqrt(np.mean((y_test - y_pred) ** 2))
r2 = r2_score(y_test, y_pred)
print(f"RMSE using only mileage predictor: {rmse:.2f}")
print(f" $R^2$  using only mileage predictor: {r2:.2f}")
```

RMSE using only mileage predictor: 14437.80

$R^2$  using only mileage predictor: 0.29

Let's visualize the tree structure

```
# Plot the tree
plt.figure(figsize=(18, 6))
plot_tree(reg_tree, feature_names=['mileage'], filled=True, rounded=True)
plt.title("Regression Tree Using Mileage")
plt.show()
```



Let's visualize how mileage is used in the decision tree below:

```
# Create evenly spaced mileage values within the range of training data
Xtest = np.linspace(X_train_mileage['mileage'].min(), X_train_mileage['mileage'].max(), 100)

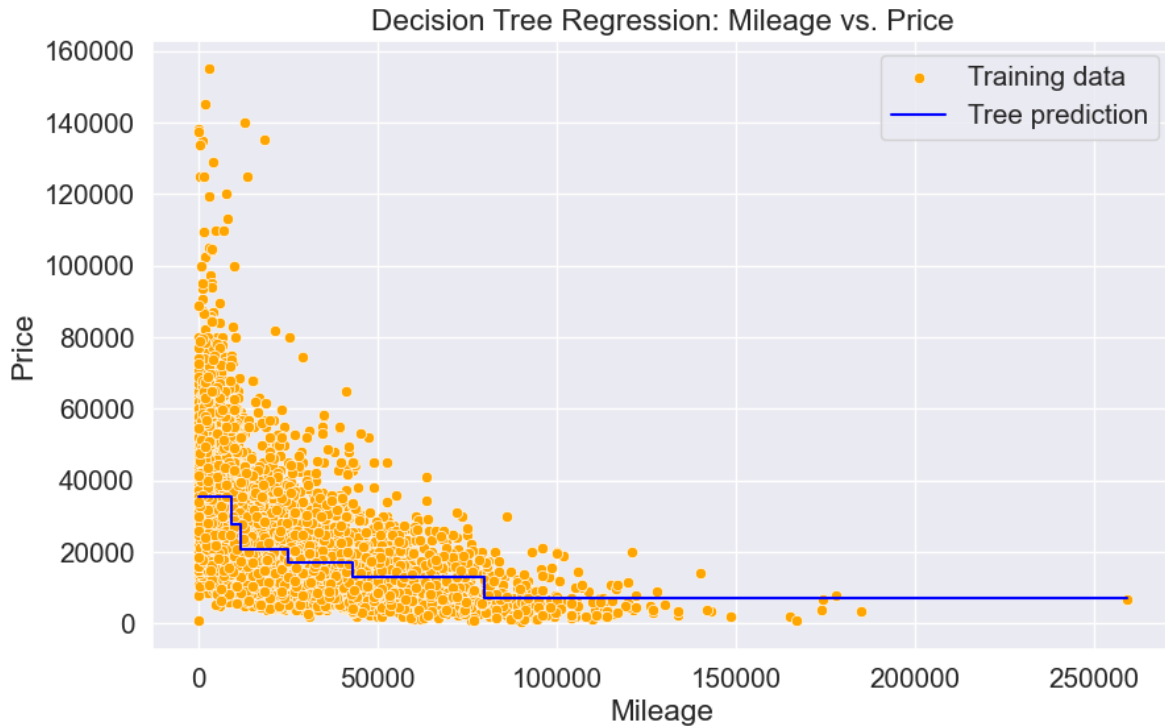
# Convert Xtest to a DataFrame with the correct column name
Xtest_df = pd.DataFrame(Xtest, columns=['mileage'])

# Predict using the DataFrame instead of NumPy array
ytest_pred = reg_tree.predict(Xtest_df)

plt.figure(figsize=(10, 6))
sns.scatterplot(x=X_train_mileage['mileage'], y=y_train, color='orange', label='Training data')

# Step plot to reflect piecewise constant predictions
plt.step(Xtest_df['mileage'], ytest_pred, color='blue', label='Tree prediction', where='mid')

plt.xlabel("Mileage")
plt.ylabel("Price")
plt.title("Decision Tree Regression: Mileage vs. Price")
plt.legend()
plt.show()
```



All cars falling within the same terminal node have the same predicted price, which is seen as flat line segments in the above model curve.

#### 4.2.2 Using mileage and brand as predictors

```
X_train.head()
```

	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize
216	vw	Scirocco	2016	Manual	41167	Diesel	20	55.2654	2.0
4381	merc	CLS Class	2018	Semi-Auto	12078	Diesel	145	47.7624	2.9
6891	hyundi	Santa Fe	2019	Automatic	623	Diesel	145	43.0887	2.2
421	hyundi	IX35	2014	Manual	37095	Diesel	145	53.4862	1.7
505	ford	Edge	2016	Semi-Auto	15727	Diesel	160	49.0741	2.0

```
# Select features and target
X_train_tree = X_train[['mileage', 'brand']]

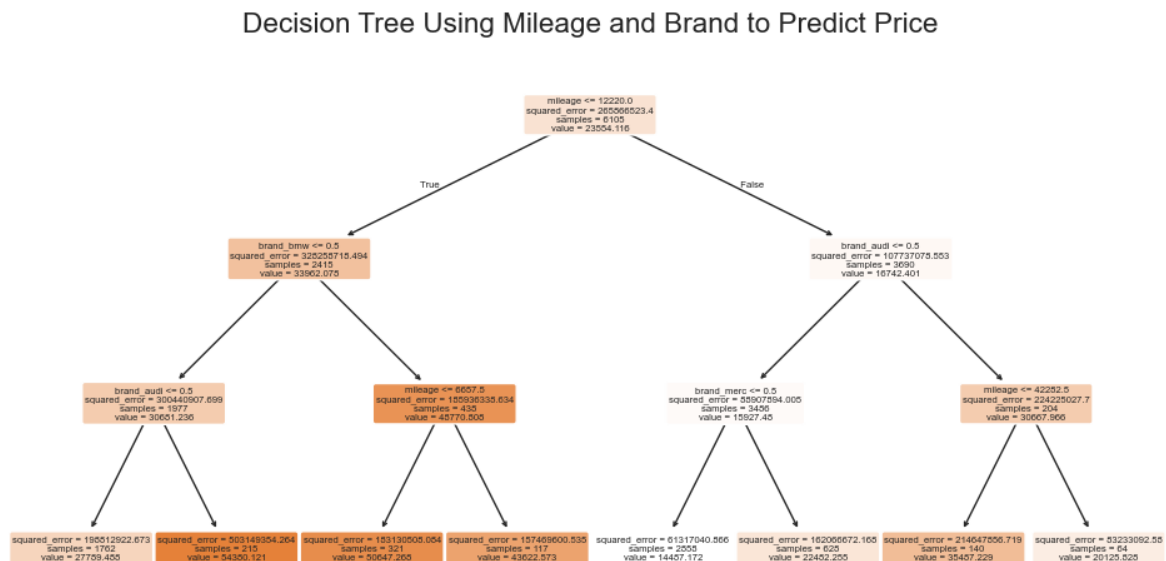
X_test_tree = X_test[['mileage', 'brand']]
```

```
# One-hot encode the categorical variable 'brand'
X_train_tree_encoded = pd.get_dummies(X_train_tree, columns=['brand'])
X_test_tree_encoded = pd.get_dummies(X_test_tree, columns=['brand'])
```

```
model = DecisionTreeRegressor(max_depth=3, random_state=42)
model.fit(X_train_tree_encoded, y_train)
```

```
DecisionTreeRegressor(max_depth=3, random_state=42)
```

```
plt.figure(figsize=(12, 6))
plot_tree(model, feature_names=X_train_tree_encoded.columns, filled=True, rounded=True)
plt.title("Decision Tree Using Mileage and Brand to Predict Price")
plt.show()
```



```
# Predict the target variable using the model
y_pred_tree = model.predict(X_test_tree_encoded)

# Calculate the RMSE and R2 score
rmse_tree = np.sqrt(np.mean((y_test - y_pred_tree) ** 2))
r2_tree = r2_score(y_test, y_pred_tree)
print(f"RMSE using mileage and brand predictor: {rmse_tree:.2f}")
print(f"R2 using mileage and brand predictor: {r2_tree:.2f}")
```

```
# Compare the performance of the two models
print(f"RMSE using only mileage predictor: {rmse:.2f}")
print(f"RMSE using mileage and brand predictor: {rmse_tree:.2f}")
# The RMSE using mileage and brand predictor is lower than using only mileage predictor.
# This indicates that adding the brand feature improves the model's performance
```

```
RMSE using mileage and brand predictor: 12531.44
R2 using mileage and brand predictor: 0.46
RMSE using only mileage predictor: 14437.80
RMSE using mileage and brand predictor: 12531.44
```

### 4.2.3 Using all predictors

Now that we've explored a single predictor (**mileage**) and added a second predictor (**brand**), let's take it a step further and use **all available features** to build a more robust model.

We'll construct a **pipeline** that handles necessary preprocessing steps (e.g., categorical encoding) and fits a **Decision Tree Regressor** in a streamlined and reproducible way.

```
# extract the categorical columns and put them in a list
categorical_feature = X.select_dtypes(include=['object']).columns.tolist()

# extract the numerical columns and put them in a list
numerical_feature = X.select_dtypes(include=['int64', 'float64']).columns.tolist()

# Create a ColumnTransformer to handle encoding
preprocessor = ColumnTransformer(
    transformers=[
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_feature)
    ],
    remainder='passthrough', # Keep numerical feature (mileage) unchanged
    force_int_remainder_cols=False
)

# Create the pipeline
pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('regressor', DecisionTreeRegressor(max_depth=4, random_state=42))
])
```

```
# Usage:
pipeline.fit(X_train, y_train)
y_pred = pipeline.predict(X_test)

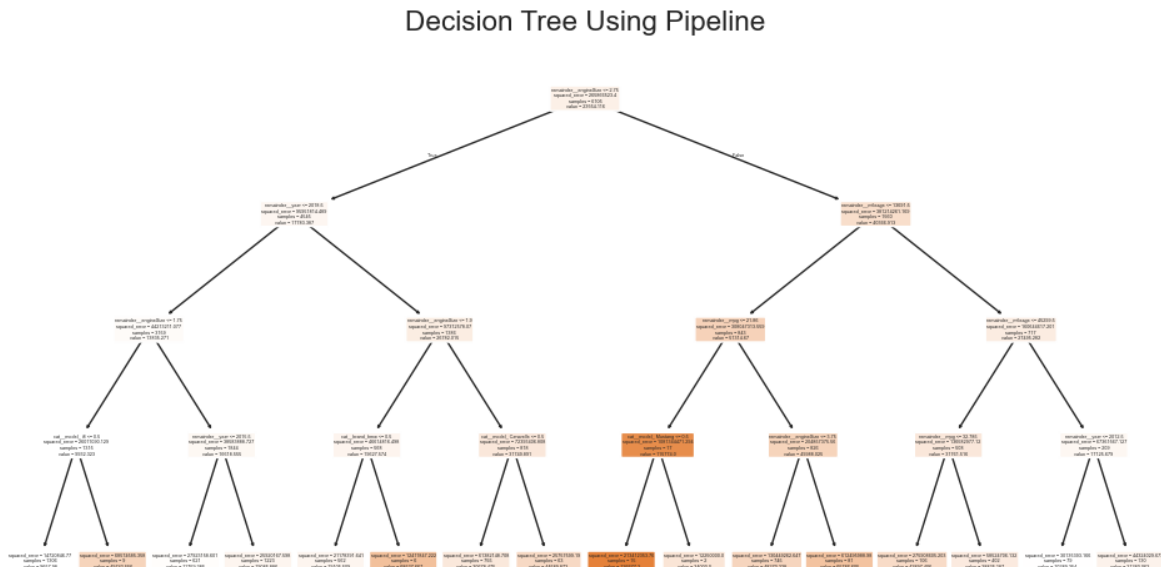
# Calculate the RMSE and R2 score
rmse_pipeline = np.sqrt(np.mean((y_test - y_pred) ** 2))
r2_pipeline = r2_score(y_test, y_pred)
print(f"RMSE using pipeline: {rmse_pipeline:.2f}")
print(f"R2 using pipeline: {r2_pipeline:.2f}")
```

RMSE using pipeline: 8186.34

R<sup>2</sup> using pipeline: 0.77

```
# let's visualize the decision tree
```

```
plt.figure(figsize=(12, 6))
plot_tree(pipeline.named_steps['regressor'], feature_names=pipeline.named_steps['preprocessor'].feature_names_in_)
plt.title("Decision Tree Using Pipeline")
plt.show()
```



## 4.3 Key Hyperparameters in Decision Tree

In regression trees, **model complexity** is controlled by hyperparameters, tuning them is crucial for balancing **underfitting** and **overfitting**.

### 4.3.1 Underfitting

- The model is **too simple** to capture patterns in the data.
- High bias, low variance.
- Often caused by:
  - Shallow trees (`max_depth` is too small)
  - Too strict constraints (`min_samples_split` or `min_samples_leaf` is too high)

### 4.3.2 Overfitting

- The model is **too complex** and learns noise from the training data.
- Low bias, high variance.
- Often caused by:
  - Deep trees with many splits
  - Very small `min_samples_leaf` or `min_samples_split`

Below are the most commonly used hyperparameters:

#### 4.3.3 `max_depth`

- Controls the **maximum depth** of the tree.
- If `None`, the tree will expand until all leaves are pure or contain fewer than `min_samples_split` samples.
- Controls overfitting (deep trees → overfit, shallow trees → underfit)
- Typical values: 3 to 20 (start with lower values).

#### 4.3.4 `min_samples_split`

- The **minimum number of samples** required to split an internal node.
- higher values → simpler trees (reducing overfitting)

### 4.3.5 min\_samples\_leaf

- The minimum number of samples required to be at a leaf node.
- Setting this to a higher number can smooth the model by reducing variance.

### 4.3.6 max\_features

- Number of features to consider when looking for the best split.
- Can be set to:
  - "auto" or None: use all features
  - "sqrt": use the square root of the number of features
  - "log2": use log base 2

```
# Define your parameter grid with pipeline step prefix
param_grid = {
    'regressor__max_depth': [3, 5, 7, 10, None],
    'regressor__min_samples_split': [2, 5, 10],
    'regressor__min_samples_leaf': [1, 2, 4],
    'regressor__max_features': ['sqrt', None]
}

# Create custom scorer for RMSE
rmse_scorer = make_scorer(lambda y_true, y_pred: root_mean_squared_error(y_true, y_pred),
                           greater_is_better=False)

# Create GridSearchCV object
grid_search = GridSearchCV(
    estimator=pipeline,
    param_grid=param_grid,
    scoring={
        'RMSE': rmse_scorer,
        'R2': 'r2'
    },
    refit='R2',
    cv=5,
    n_jobs=-1,
    verbose=1
)

# Fit the grid search to training data
grid_search.fit(X_train, y_train)
```



Fitting 5 folds for each of 270 candidates, totalling 1350 fits

```
GridSearchCV(cv=5,
             estimator=Pipeline(steps=[('preprocessor',
                                       ColumnTransformer(force_int_remainder_cols=False,
                                                         remainder='passthrough',
                                                         transformers=[('cat',
                                                                       OneHotEncoder(handle_unknown='ignore',
                                                                       sparse_output=False),
                                                                       ['brand',
                                                                       'model',
                                                                       'transmission',
                                                                       'fuelType'])])),
                                     ('regressor',
                                      DecisionTreeRegressor(max_depth=4,
                                                             random_state=42))])),
             n_jobs=-1,
             param_grid={'regressor__ccp_alpha': [0.001, 0.01, 0.1],
                         'regressor__max_depth': [3, 5, 7, 10, None],
                         'regressor__max_features': ['sqrt', None],
                         'regressor__min_samples_leaf': [1, 2, 4],
                         'regressor__min_samples_split': [2, 5, 10]},
             refit='R2',
             scoring={'R2': 'r2',
                     'RMSE': make_scorer(<lambda>, greater_is_better=False, response_method='raw')},
             verbose=1)
```

The `GridSearchCV` setup evaluates both **RMSE** and **R<sup>2</sup>** during cross-validation.

- **R<sup>2</sup>** is used to **select the best model** and is also used to **refit** the model on the entire training set.
- **RMSE** is computed during the process for evaluation purposes, but it is **not used to determine** the best model.

This allows for more comprehensive model assessment while still optimizing based on a single selected metric.

```
# Get best estimator and predictions
best_model = grid_search.best_estimator_
y_pred_tuned = best_model.predict(X_test)

# Calculate metrics for tuned model
rmse_tuned = root_mean_squared_error(y_test, y_pred_tuned)
r2_tuned = r2_score(y_test, y_pred_tuned)
```

```

print("\n=== Best Parameters ===")
print(grid_search.best_params_)
print("\n=== Tuned Model Performance ===")
print(f"RMSE (Tuned): {rmse_tuned:.2f}")
print(f"R2 (Tuned): {r2_tuned:.2f}")
print(f"Improvement in R2: {(r2_tuned - r2_pipeline):.2%}")

```

```

=== Best Parameters ===
{'regressor__ccp_alpha': 0.001, 'regressor__max_depth': None, 'regressor__max_features': Non

=== Tuned Model Performance ===
RMSE (Tuned): 4726.17
R2 (Tuned): 0.92
Improvement in R2: 15.23%

```

```

print("\n=== Best Parameters ===")
print(grid_search.best_params_)
print("\n=== Tuned Model Performance ===")
print(f"RMSE (Tuned): {rmse_tuned:.2f}")
print(f"R2 (Tuned): {r2_tuned:.2f}")
print(f"Improvement in R2: {(r2_tuned - r2_pipeline):.2%}")

```

```

=== Best Parameters ===
{'regressor__ccp_alpha': 0.001, 'regressor__max_depth': None, 'regressor__max_features': Non

=== Tuned Model Performance ===
RMSE (Tuned): 4726.17
R2 (Tuned): 0.92
Improvement in R2: 15.23%

```

GridSearchCV improves the r squared from 0.77 to 0.92, increased by 15.23%, Let us visualize the mean squared error based on the hyperparameter values. We'll use the cross validation results stored in the cv\_results\_ attribute of the GridSearchCV fit() object.

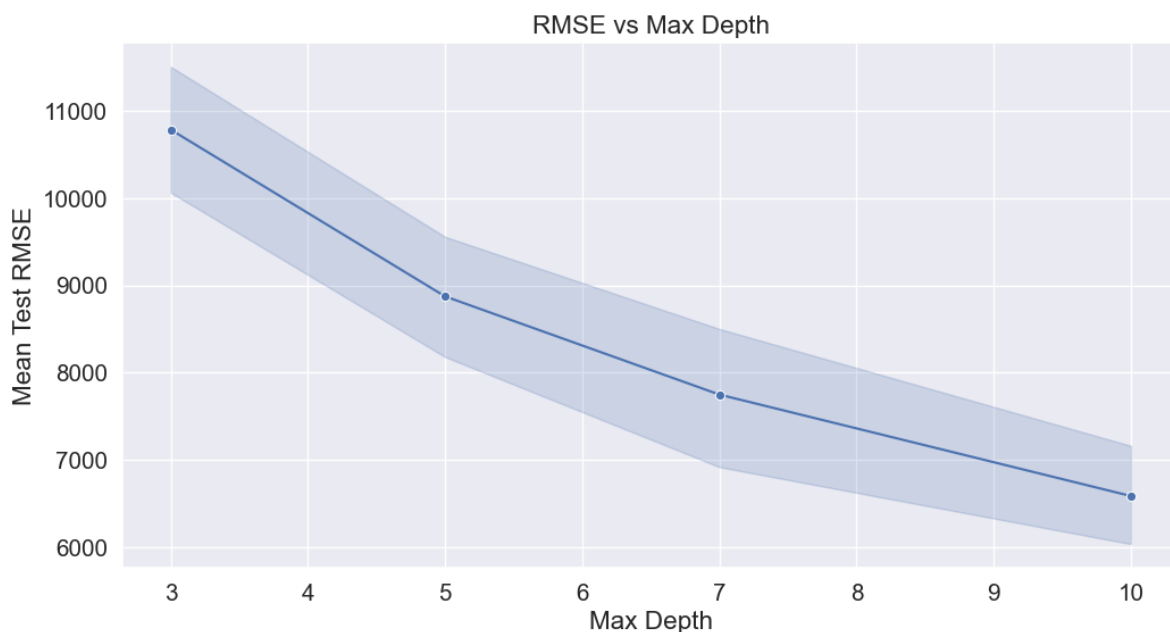
```

#Detailed results of k-fold cross validation
cv_results = pd.DataFrame(grid_search.cv_results_)
cv_results.head()

```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_regressor__ccp_alpha	p
0	0.075184	0.013049	0.009999	0.001052	0.001	3
1	0.078505	0.014115	0.010804	0.001051	0.001	3
2	0.077370	0.011081	0.010873	0.000946	0.001	3
3	0.036274	0.036182	0.010617	0.000345	0.001	3
4	0.018702	0.003564	0.012120	0.004090	0.001	3

```
# Plotting the RMSE for different max_depth values
plt.figure(figsize=(12, 6))
sns.lineplot(data=cv_results, x='param_regressor__max_depth', y=np.abs(cv_results['mean_test_
plt.xlabel('Max Depth')
plt.ylabel('Mean Test RMSE')
plt.title('RMSE vs Max Depth');
```



#### 4.3.7 Output feature importance

```
# Get feature importances and names
feature_importances = best_model.named_steps['regressor'].feature_importances_
feature_names = best_model.named_steps['preprocessor'].get_feature_names_out()
```

```

# Create DataFrame and select top 10
feature_importance_df = (
    pd.DataFrame({'Feature': feature_names, 'Importance': feature_importances})
    .sort_values(by='Importance', ascending=False)
    .head(10) # Keep only top 10 features
)

# Print top 10 features
print("=== Top 10 Feature Importances ===")
print(feature_importance_df)

# Plot top 10 features
plt.figure(figsize=(12, 6))
sns.barplot(data=feature_importance_df, x='Importance', y='Feature')
plt.title('Top 10 Feature Importances from Decision Tree')
plt.xlabel('Importance')
plt.ylabel('Feature')
plt.tight_layout()
plt.show()

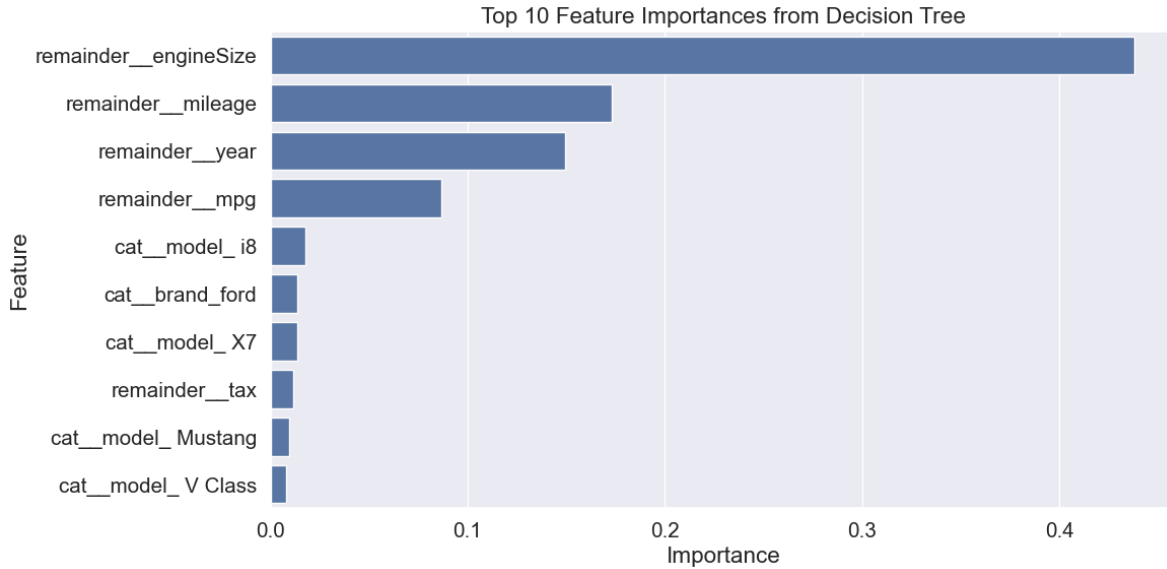
```

```

=== Top 10 Feature Importances ===

```

	Feature	Importance
112	remainder__engineSize	0.437921
109	remainder__mileage	0.173215
108	remainder__year	0.149303
111	remainder__mpg	0.086922
98	cat__model_ i8	0.017971
2	cat__brand_ford	0.013837
92	cat__model_ X7	0.013477
110	remainder__tax	0.011502
60	cat__model_ Mustang	0.009517
86	cat__model_ V Class	0.008147



## 4.4 Cost-Complexity Pruning (`ccp_alpha`)

Cost-complexity pruning is a post-pruning technique used to reduce the size of a decision tree by removing sections that provide little to no improvement in prediction accuracy. It helps prevent **overfitting** and improves **model generalization**.

### 4.4.1 Key Idea

Each subtree in a decision tree has an associated **cost-complexity score**:

$$R_{\alpha}(T) = R(T) + \alpha |T|$$

- $R(T)$ : Total training error of the tree ( $T$ )
- $|T|$ : Number of leaf nodes in the tree
- $\alpha$  (**`ccp_alpha`**): Complexity parameter that penalizes tree size

As  $\alpha$  increases, the tree is **pruned more aggressively**.

### 4.4.2 Parameter: `ccp_alpha` in scikit-learn

- Available in `DecisionTreeRegressor` and `DecisionTreeClassifier`
- Default: `ccp_alpha = 0.0` (no pruning)
- Increasing `ccp_alpha` encourages simpler trees by penalizing extra leaf nodes

```

from sklearn.model_selection import GridSearchCV
from sklearn.metrics import make_scorer, mean_squared_error
import numpy as np

# Define your parameter grid with pipeline step prefix
param_grid = {
    'regressor__ccp_alpha': [0.0, 0.001, 0.01, 0.1]
}

# Create custom scorer for RMSE
rmse_scorer = make_scorer(lambda y_true, y_pred: np.sqrt(mean_squared_error(y_true, y_pred)),
                           greater_is_better=False)

# Create GridSearchCV object
grid_search_ccp = GridSearchCV(
    estimator=pipeline,
    param_grid=param_grid,
    scoring={
        'RMSE': rmse_scorer,
        'R2': 'r2'
    },
    refit='R2',
    cv=5,
    n_jobs=-1,
    verbose=1
)

# Fit the grid search to training data
grid_search_ccp.fit(X_train, y_train)

```

Fitting 5 folds for each of 4 candidates, totalling 20 fits

```

GridSearchCV(cv=5,
             estimator=Pipeline(steps=[('preprocessor',
                                       ColumnTransformer(force_int_remainder_cols=False,
                                                         remainder='passthrough',
                                                         transformers=[('cat',
                                                                       OneHotEncoder(handle_unknown='ignore'),
                                                                       ['brand',
                                                                        'model',
                                                                        'transmission',

```

```

('regressor',
 DecisionTreeRegressor(max_depth=4,
                        random_state=42))]),
n_jobs=-1,
param_grid={'regressor__ccp_alpha': [0.0, 0.001, 0.01, 0.1]},
refit='R2',
scoring={'R2': 'r2',
         'RMSE': make_scorer(<lambda>, greater_is_better=False, response_method=
verbose=1)

```

```

encoder = OneHotEncoder(handle_unknown='ignore', sparse_output=False)

X_train_encoded = encoder.fit_transform(X_train[categorical_feature])
X_test_encoded = encoder.transform(X_test[categorical_feature])

# Convert the encoded features back to DataFrame
X_train_encoded_df = pd.DataFrame(X_train_encoded, columns=encoder.get_feature_names_out(categorical_feature))
X_test_encoded_df = pd.DataFrame(X_test_encoded, columns=encoder.get_feature_names_out(categorical_feature))

# Concatenate the encoded features with the original numerical features
X_train_final = pd.concat([X_train_encoded_df, X_train[numerical_feature].reset_index(drop=True)], axis=1)
X_test_final = pd.concat([X_test_encoded_df, X_test[numerical_feature].reset_index(drop=True)], axis=1)

# Check the final shape of the training and testing sets
print("Training set shape:", X_train_final.shape)
print("Testing set shape:", X_test_final.shape)
# Check the first few rows of the final training set
X_train_final.head()
# Check the first few rows of the final testing set
X_test_final.head()

```

Training set shape: (6105, 113)

Testing set shape: (1527, 113)

	brand_audi	brand_bmw	brand_ford	brand_hyundi	brand_merc	brand_skoda	brand_toyota
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0	1.0	0.0	0.0
2	0.0	0.0	0.0	1.0	0.0	0.0	0.0
3	0.0	0.0	0.0	1.0	0.0	0.0	0.0
4	0.0	0.0	1.0	0.0	0.0	0.0	0.0

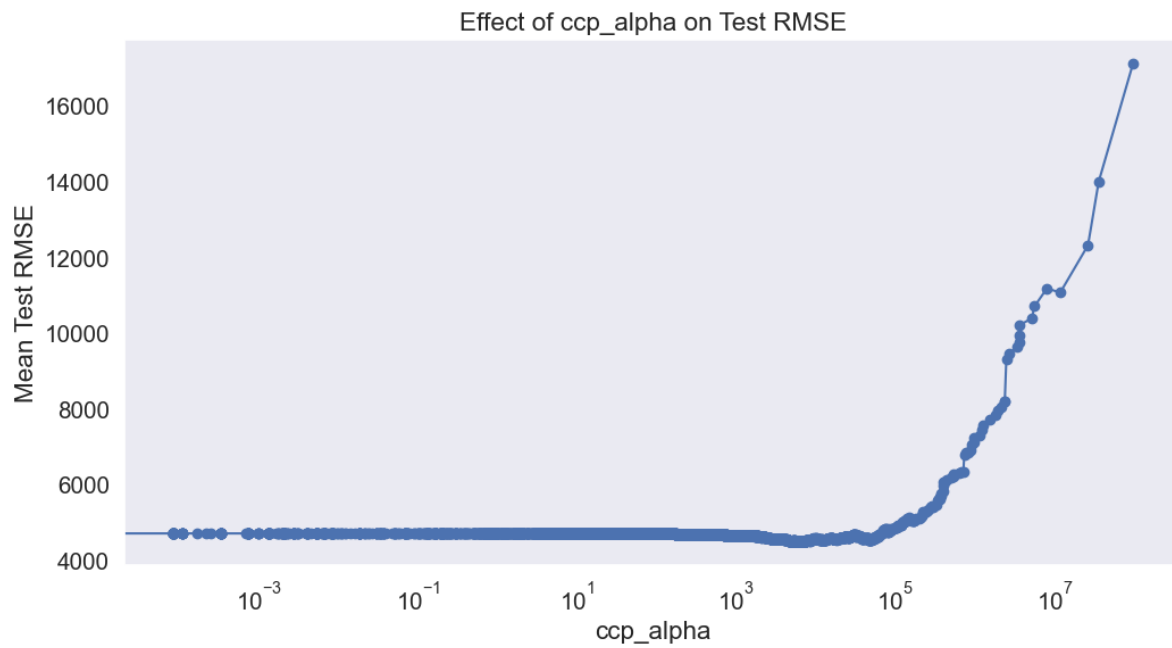
```

model = DecisionTreeRegressor(random_state = 1)#model without any restrictions
path= model.cost_complexity_pruning_path(X_train_final,y_train)# Compute the pruning path du

# Extract the effective alphas and the corresponding performance metrics
ccp_alphas = path.ccp_alphas
impurities = path.impurities
# Create a DataFrame to store the results
ccp_results = pd.DataFrame({'ccp_alpha': ccp_alphas, 'impurity': impurities})
# Fit the model for each alpha value and calculate the mean test score
mean_test_scores = []
for alpha in ccp_alphas:
    model = DecisionTreeRegressor(random_state=1, ccp_alpha=alpha)
    model.fit(X_train_final, y_train)
    y_pred = model.predict(X_test_final)
    mean_test_scores.append(np.sqrt(mean_squared_error(y_test, y_pred)))
# Add the mean test scores to the DataFrame
ccp_results['mean_test_score'] = mean_test_scores
# Plot the results
plt.figure(figsize=(12, 6))
plt.plot(ccp_results['ccp_alpha'], ccp_results['mean_test_score'], marker='o')
plt.xlabel('ccp_alpha')
plt.ylabel('Mean Test RMSE')
plt.title('Effect of ccp_alpha on Test RMSE')
plt.xscale('log')
plt.grid()

```





## 5 Classification trees

*Read section 8.1.2 of the book before using these notes.*

*Note that in this course, lecture notes are not sufficient, you must read the book for better understanding. Lecture notes are just implementing the concepts of the book on a dataset, but not explaining the concepts elaborately.*

Import libraries

```
# %load ../standard_import.txt
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from IPython.display import Image

from sklearn.model_selection import train_test_split, cross_val_score
from six import StringIO
from sklearn.tree import export_graphviz, DecisionTreeClassifier, plot_tree
from sklearn.metrics import confusion_matrix, accuracy_score

%matplotlib inline

# load the dataset
heart_df = pd.read_csv('datasets/heart_disease_classification.csv')
print(heart_df.shape)
heart_df.head()
```

(303, 14)

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

```
# print out target distribution
heart_df.target.value_counts()
```

```
target
1    165
0    138
Name: count, dtype: int64
```

```
# split the x and y data
X = heart_df.drop(columns=['target'])
y = heart_df.target

# split the data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

## 5.1 Building a Classification Tree

We will build a classification tree to predict whether a person has heart disease, using the default parameters of the decision tree classifier.

```
# create a decision tree classifier
tree = DecisionTreeClassifier(random_state=42)

# fit the model to the training data
tree.fit(X_train, y_train)

# make predictions on the test data
y_pred = tree.predict(X_test)

# calculate the accuracy of the model
accuracy = accuracy_score(y_test, y_pred)
print(f'Test Accuracy: {accuracy:.2f}')

# train accuracy
```



```
tree_depth = tree.get_depth()
print(f'Depth of the tree: {tree_depth}')
```

Number of leaves: 41  
Depth of the tree: 9

Clearly, the model is overfitting, as indicated by a training accuracy of 100% and a much lower test accuracy of 75%.

Next, we will explore different strategies to address and reduce overfitting.

## 5.2 Pre-pruning: Hyperparameters Tuning

Maximum depth of tree (**max\_depth**) - Used to control over-fitting as higher depth will allow model to learn relations very specific to a particular sample.

Minimum samples for a node split (**min\_samples\_split**) - Defines the minimum number of samples (or observations) which are required in a node to be considered for splitting. - Used to control over-fitting. Higher values prevent a model from learning relations which might be highly specific to the particular sample selected for a tree.

Minimum samples for a terminal node (**min\_samples\_leaf**) - Defines the minimum samples (or observations) required in a terminal node or leaf. - Used to control over-fitting similar to **min\_samples\_split**. - Generally lower values should be chosen for imbalanced class problems because the regions in which the minority class will be in majority will be very small.

Maximum number of terminal nodes (**max\_leaf\_nodes**) - The maximum number of terminal nodes or leaves in a tree.

```
# hyperparameter tuning

from sklearn.model_selection import GridSearchCV

# define the parameter grid
param_grid = {
    'max_depth': list(range(1, 9)) + [None],
    'min_samples_split': [2, 5, 10, 15, 20],
    'min_samples_leaf': [1, 2, 4]
}

# create a grid search object
grid_search = GridSearchCV(estimator=tree, param_grid=param_grid, cv=5, n_jobs=-1, verbose=2)
```

```

# fit the grid search to the training data
grid_search.fit(X_train, y_train)
# print the best parameters
print("Best parameters found: ", grid_search.best_params_)

# print the best score
print("Best score: ", grid_search.best_score_)
# get the best estimator
best_tree = grid_search.best_estimator_

# make predictions on the test data with the best estimator
y_pred_best = best_tree.predict(X_test)
# calculate the accuracy of the best estimator
best_accuracy = accuracy_score(y_test, y_pred_best)
print(f'Best Test Accuracy: {best_accuracy:.2f}')

```

Fitting 5 folds for each of 135 candidates, totalling 675 fits  
 Best parameters found: {'max\_depth': 6, 'min\_samples\_leaf': 2, 'min\_samples\_split': 10}  
 Best score: 0.7687074829931972  
 Best Test Accuracy: 0.85

```

# print out the best tree depth and number of leaves
best_num_leaves = best_tree.get_n_leaves()
print(f'Best Number of leaves: {best_num_leaves}')

best_tree_depth = best_tree.get_depth()
print(f'Best Depth of the tree: {best_tree_depth}')

```

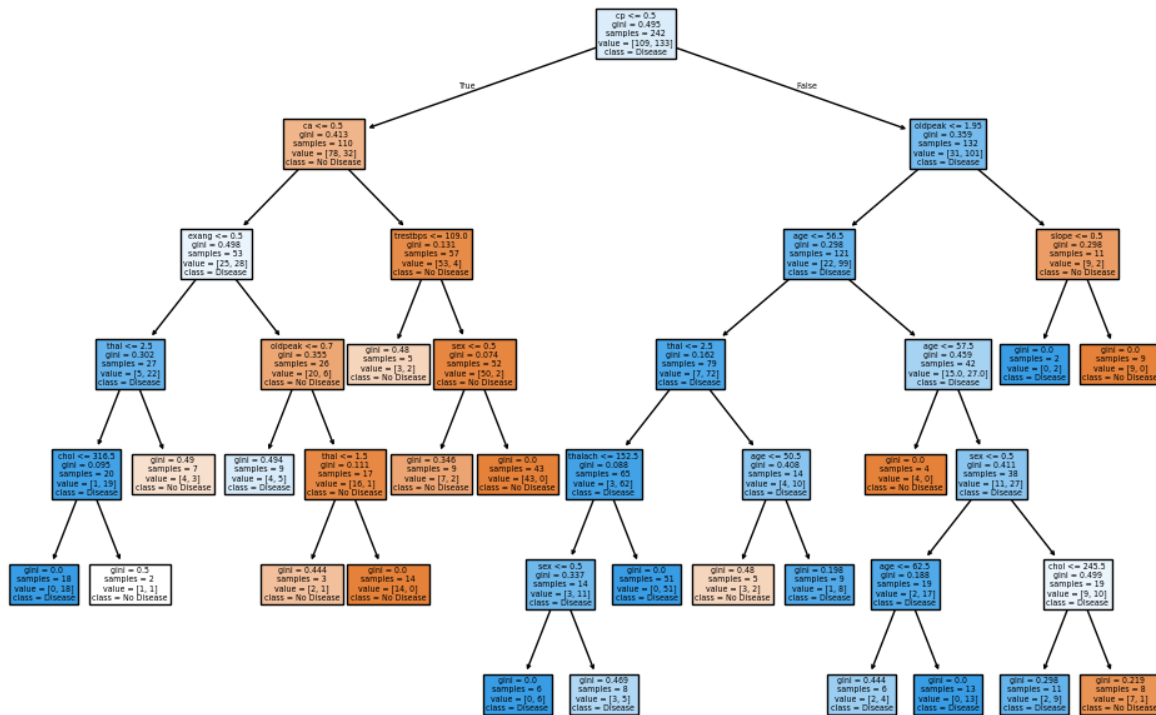
Best Number of leaves: 21  
 Best Depth of the tree: 6

```

# plot the best decision tree
plt.figure(figsize=(12, 8))
plot_tree(best_tree, filled=True, feature_names=X.columns, class_names=['No Disease', 'Disease'])
plt.title('Best Decision Tree for Heart Disease Classification')
plt.show()

```

Best Decision Tree for Heart Disease Classification



### 5.2.1 Gini or entropy

```
# define the parameter grid
param_grid_metric = {
    'max_depth': list(range(1, 9)) + [None],
    'min_samples_split': [2, 5, 10, 15, 20],
    'min_samples_leaf': [1, 2, 4],
    # adding the criterion parameter to the grid search
    'criterion': ['gini', 'entropy']
}

# create a grid search object
grid_search_metric = GridSearchCV(estimator=tree, param_grid=param_grid_metric, cv=5, n_jobs=
# fit the grid search to the training data
grid_search_metric.fit(X_train, y_train)
# print the best parameters
print("Best parameters found: ", grid_search_metric.best_params_)
```

```

# print the best score
print("Best score: ", grid_search_metric.best_score_)
# get the best estimator
best_tree = grid_search_metric.best_estimator_

# make predictions on the test data with the best estimator
y_pred_best = best_tree.predict(X_test)
# calculate the accuracy of the best estimator
best_accuracy = accuracy_score(y_test, y_pred_best)
print(f'Best Test Accuracy: {best_accuracy:.2f}')

```

Fitting 5 folds for each of 270 candidates, totalling 1350 fits

Best parameters found: {'criterion': 'entropy', 'max\_depth': 4, 'min\_samples\_leaf': 1, 'min

Best score: 0.7811224489795918

Best Test Accuracy: 0.85

```

# print out the best tree depth and number of leaves
best_num_leaves = best_tree.get_n_leaves()
print(f'Best Number of leaves: {best_num_leaves}')

best_tree_depth = best_tree.get_depth()
print(f'Best Depth of the tree: {best_tree_depth}')

```

Best Number of leaves: 11

Best Depth of the tree: 4

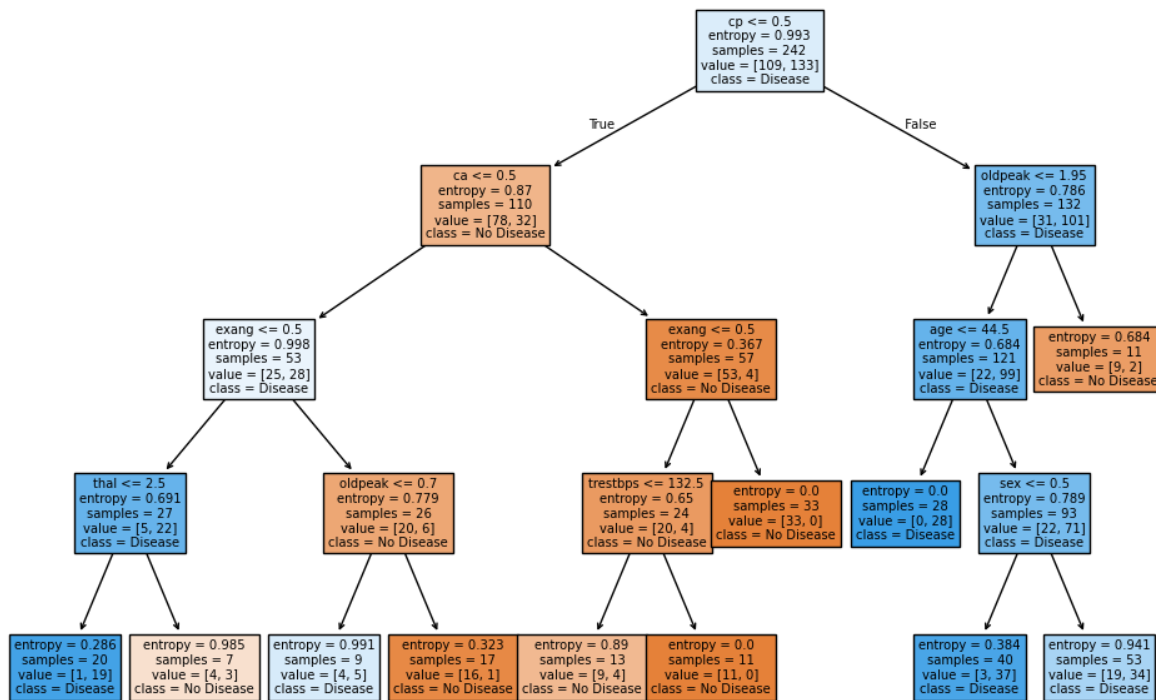
```

# plot the best decision tree
plt.figure(figsize=(12, 8))
plot_tree(best_tree, filled=True, feature_names=X.columns, class_names=['No Disease', 'Disease'])
plt.title('Best Decision Tree for Heart Disease Classification')
plt.show()

```



Best Decision Tree for Heart Disease Classification



Both criteria aim to minimize impurity in splits, in practice, they often lead to comparable performance in decision trees, even though their mathematical formulations differ

- **Gini:** Faster to compute (no logarithms) and often used for large datasets. so it is a good default. May produce slightly more complex trees.
- **Entropy:** Slower but aligns with information theory. Prefers splits that balance node sizes, leading to more interpretable trees.

### 5.3 Post-pruning: Cost complexity pruning

Post-pruning, on the other hand, allows the decision tree to grow to its full extent and then prunes it back to reduce complexity. This approach first builds a complete tree and then removes or collapses branches that don't significantly contribute to the model's performance. One common post-pruning technique is called Cost-Complexity Pruning.

### 5.3.1 step 1: calculate the cost complexity pruning path

```
path = tree.cost_complexity_pruning_path(X_train, y_train)
ccp_alphas, impurities = path.ccp_alphas, path.impurities
```

### 5.3.2 step 2: Create trees with different ccp\_alpha values and evaluate their performance

```
# We'll skip the last alpha which would produce a single-node tree
alphas = ccp_alphas[:-1]

# Create empty lists to store the results
train_scores = []
test_scores = []
cv_scores = []
node_counts = []

# For each alpha value, fit a tree and evaluate
for alpha in alphas:
    # Create and train the model
    clf = DecisionTreeClassifier(ccp_alpha=alpha, random_state=42)
    clf.fit(X_train, y_train)

    # Record scores
    train_scores.append(accuracy_score(y_train, clf.predict(X_train)))
    test_scores.append(accuracy_score(y_test, clf.predict(X_test)))

    # Cross-validation score for robustness
    cv_score = cross_val_score(clf, X_train, y_train, cv=5).mean()
    cv_scores.append(cv_score)

    # Record tree complexity
    node_counts.append(clf.tree_.node_count)
```

### 5.3.3 Step 3: Visualize the results

```

# Step 3: Visualize the results
fig, ax = plt.subplots(2, 2, figsize=(15, 10))

# Plot accuracy vs alpha
ax[0, 0].plot(alphas, train_scores, marker='o', label='Train')
ax[0, 0].plot(alphas, test_scores, marker='o', label='Test')
ax[0, 0].plot(alphas, cv_scores, marker='o', label='Cross-validation')
ax[0, 0].set_xlabel('ccp_alpha')
ax[0, 0].set_ylabel('Accuracy')
ax[0, 0].set_title('Accuracy vs. ccp_alpha')
ax[0, 0].legend()
ax[0, 0].grid(True)

# Plot number of nodes vs alpha
ax[0, 1].plot(alphas, node_counts, marker='o')
ax[0, 1].set_xlabel('ccp_alpha')
ax[0, 1].set_ylabel('Number of nodes')
ax[0, 1].set_title('Tree complexity vs. ccp_alpha')
ax[0, 1].grid(True)

# Log scale for better visualization of small alpha values
ax[1, 0].plot(alphas, train_scores, marker='o', label='Train')
ax[1, 0].plot(alphas, test_scores, marker='o', label='Test')
ax[1, 0].plot(alphas, cv_scores, marker='o', label='Cross-validation')
ax[1, 0].set_xlabel('ccp_alpha (log scale)')
ax[1, 0].set_ylabel('Accuracy')
ax[1, 0].set_title('Accuracy vs. ccp_alpha (log scale)')
ax[1, 0].set_xscale('log')
ax[1, 0].legend()
ax[1, 0].grid(True)

# Find best alpha based on test score
best_test_idx = np.argmax(test_scores)
best_test_alpha = alphas[best_test_idx]
best_test_acc = test_scores[best_test_idx]

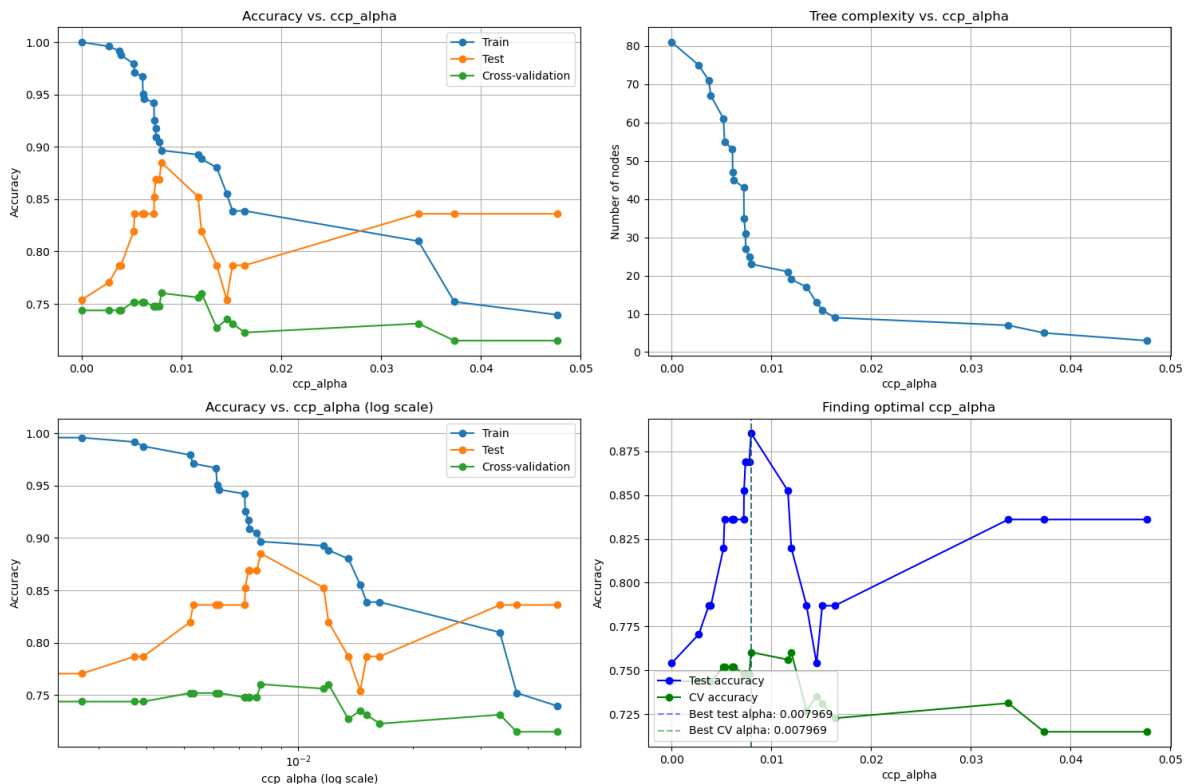
# Find best alpha based on CV score (more robust)
best_cv_idx = np.argmax(cv_scores)
best_cv_alpha = alphas[best_cv_idx]
best_cv_acc = cv_scores[best_cv_idx]

```

```
# Plot highlighting best points
ax[1, 1].plot(alphas, test_scores, 'b-', marker='o', label='Test accuracy')
ax[1, 1].plot(alphas, cv_scores, 'g-', marker='o', label='CV accuracy')
ax[1, 1].axvline(x=best_test_alpha, color='blue', linestyle='--', alpha=0.5,
                 label=f'Best test alpha: {best_test_alpha:.6f}')
ax[1, 1].axvline(x=best_cv_alpha, color='green', linestyle='--', alpha=0.5,
                 label=f'Best CV alpha: {best_cv_alpha:.6f}')
ax[1, 1].set_xlabel('ccp_alpha')
ax[1, 1].set_ylabel('Accuracy')
ax[1, 1].set_title('Finding optimal ccp_alpha')
ax[1, 1].legend(loc='lower left')
ax[1, 1].grid(True)

plt.tight_layout()
plt.show()

# Print the optimal alpha values and corresponding metrics
print(f"Best alpha based on test score: {best_test_alpha:.6f} (Accuracy: {best_test_acc:.4f})
print(f"Best alpha based on CV score: {best_cv_alpha:.6f} (Accuracy: {best_cv_acc:.4f}, Nodes: {best_cv_nodes:.4f})")
```



### 5.3.4 Step 4: Create the final model with the optimal alpha

```
# Using CV-based alpha as it's more robust against overfitting
final_model = DecisionTreeClassifier(ccp_alpha=best_cv_alpha, random_state=42)
final_model.fit(X_train, y_train)

# Evaluate the final model
train_acc = accuracy_score(y_train, final_model.predict(X_train))
test_acc = accuracy_score(y_test, final_model.predict(X_test))

print(f"\nFinal model performance:")
print(f"Training accuracy: {train_acc:.4f}")
print(f"Test accuracy: {test_acc:.4f}")
print(f"Tree nodes: {final_model.tree_.node_count}")
print(f"Tree depth: {final_model.get_depth()}")
```

Best alpha based on test score: 0.007969 (Accuracy: 0.8852, Nodes: 23)

Best alpha based on CV score: 0.007969 (Accuracy: 0.7604, Nodes: 23)

Final model performance:

Training accuracy: 0.8967

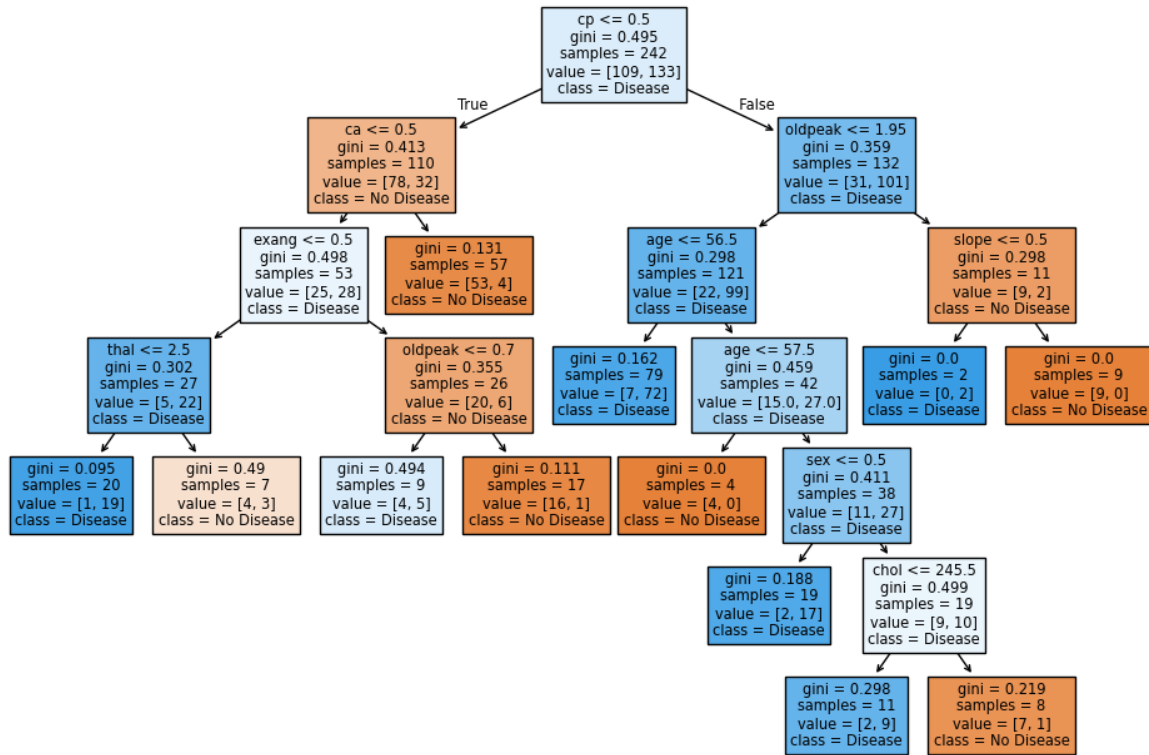
Test accuracy: 0.8852

Tree nodes: 23

Tree depth: 6

```
# plot the final decision tree
plt.figure(figsize=(12, 8))
plot_tree(final_model, filled=True, feature_names=X.columns, class_names=['No Disease', 'Disease'])
plt.title('Final Decision Tree for Heart Disease Classification')
plt.show()
```

Final Decision Tree for Heart Disease Classification



Post-pruning can potentially create more optimal trees, as it considers the entire tree structure before making pruning decisions. However, it can be more computationally expensive.

Both approaches aim to find a balance between model complexity and performance, with the goal of creating a model that generalizes well to unseen data. The choice between pre-pruning and post-pruning (or a combination of both) often depends on the specific dataset, the problem at hand, and of course, computational resources available.

## 5.4 Feature Importance in Decision Trees

Decision tree algorithms, such as Classification and Regression Trees (CART), compute **feature importance** scores based on how much each feature contributes to reducing the splitting criterion (e.g., **Gini impurity** or **entropy**).

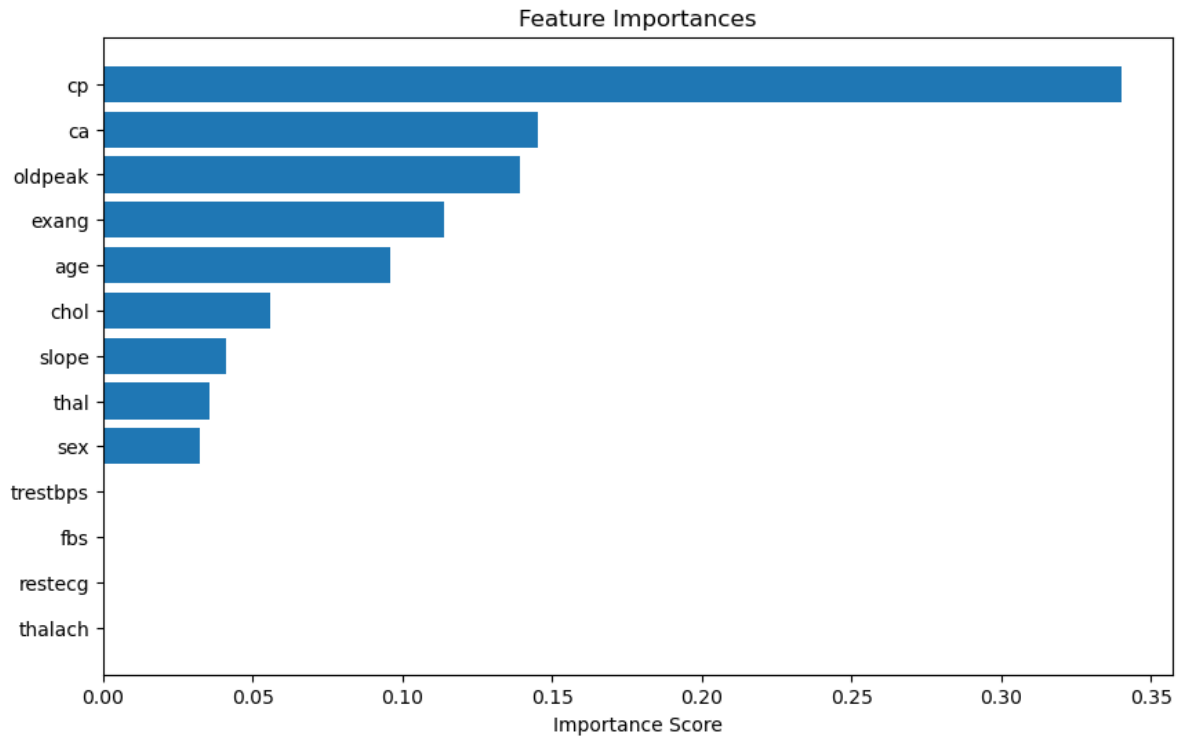
This methodology extends naturally to **ensemble models** like **Random Forests** and **Gradient Boosting**, which average feature importance across all trees in the ensemble.

Once a model is trained, the relative importance of each feature can be accessed using the `.feature_importances_` attribute. These scores indicate how valuable each feature was in constructing the decision rules that led to the model's predictions.

```
importances = final_model.feature_importances_  
  
# Create a DataFrame for easier visualization  
feature_importance_df = pd.DataFrame({  
    'Feature': X.columns,  
    'Importance': importances  
}).sort_values(by='Importance', ascending=False)  
  
# Display the DataFrame  
feature_importance_df
```

	Feature	Importance
2	cp	0.340102
11	ca	0.145273
9	oldpeak	0.139508
8	exang	0.113871
0	age	0.095885
4	chol	0.056089
10	slope	0.041242
12	thal	0.035531
1	sex	0.032499
3	trestbps	0.000000
5	fbs	0.000000
6	restecg	0.000000
7	thalach	0.000000

```
plt.figure(figsize=(10, 6))  
plt.barh(feature_importance_df['Feature'], feature_importance_df['Importance'])  
plt.xlabel("Importance Score")  
plt.title("Feature Importances")  
plt.gca().invert_yaxis() # Most important feature at the top  
plt.show()
```



#### 5.4.1 Do We Need Feature Selection or Regularization with Tree Models?

When using **tree-based models** (e.g., Decision Trees, Random Forests, Gradient Boosting) and you have a **small number of predictors**, you typically **don't need to worry much about feature selection or regularization**. Here's why:

- **Tree models inherently ignore uninformative features.** They only split on features that reduce impurity (e.g., Gini or entropy), so irrelevant features tend to receive **zero or very low importance**.
- **Ensemble methods** like Random Forest or Boosting average over many trees. Features that don't help prediction are rarely used across the ensemble.
- **Including unimportant features won't significantly hurt model performance** in small feature spaces. The model will usually ignore them during training.
- **However**, even in small feature sets, irrelevant features may slightly:
  - Increase model variance
  - Increase training time
  - Reduce interpretability



### 5.4.2 Bottom Line

If the number of predictors is small, you can typically skip feature selection and regularization when using tree-based models — they will handle irrelevant features gracefully.

## 5.5 Next Lecture

A single decision tree is highly susceptible to **overfitting**, especially on noisy or complex datasets.

In this notebook, we explored two techniques to reduce variance and improve generalization:

- **Pre-pruning** (e.g., setting `max_depth`, `min_samples_split`, etc)
- **Post-pruning** (e.g., cost-complexity pruning)

In the next lecture, we'll introduce another powerful method to combat overfitting:

**Bagging (Bootstrap Aggregating)** — an ensemble technique that builds multiple trees and averages their predictions to reduce variance and improve robustness.

## 6 Bagging

*Read section 8.2.1 of the book before using these notes.*

*Note that in this course, lecture notes are not sufficient, you must read the book for better understanding. Lecture notes are just implementing the concepts of the book on a dataset, but not explaining the concepts elaborately.*

### 6.1 Bagging: A Variance Reduction Technique

Bagging, short for **Bootstrap Aggregating**, is an effective way to reduce overfitting in decision trees. It works by training multiple decision trees—each on a different bootstrap sample of the data—and then aggregating their predictions.

Each individual decision tree is a **weak learner** and prone to overfitting, but by combining many such trees, bagging produces a **strong learner** with **lower variance** and improved generalization performance.

The number of trees to include in the ensemble is specified by the `n_estimators` hyperparameter.

```
# Importing necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(font_scale=1.35)

# import the decision tree regressor
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier, plot_tree, export_graphviz
from sklearn.ensemble import BaggingRegressor, BaggingClassifier

# split the dataset into training and testing sets
from sklearn.model_selection import train_test_split

from sklearn.model_selection import cross_val_score, GridSearchCV, cross_val_predict, KFold
```

```

from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, FunctionTransformer
from sklearn.metrics import root_mean_squared_error, r2_score, make_scorer, accuracy_score

```

## 6.2 Bagging Regression Trees

Let's revisit the same dataset used for building a single regression tree and explore whether we can further improve its performance using bagging

```

# Load the dataset
car = pd.read_csv('Datasets/car.csv')
car.head()

```

	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	vw	Beetle	2014	Manual	55457	Diesel	30	65.3266	1.6	7490
1	vauxhall	GTC	2017	Manual	15630	Petrol	145	47.2049	1.4	10998
2	merc	G Class	2012	Automatic	43000	Diesel	570	25.1172	3.0	44990
3	audi	RS5	2019	Automatic	10	Petrol	145	30.5593	2.9	51990
4	merc	X-CLASS	2018	Automatic	14000	Diesel	240	35.7168	2.3	28990

Split the predictors and target, then perform the train-test split

```

X = car.drop(columns=['price'])
y = car['price']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# extract the categorical columns and put them in a list
categorical_feature = X.select_dtypes(include=['object']).columns.tolist()

# extract the numerical columns and put them in a list
numerical_feature = X.select_dtypes(include=['int64', 'float64']).columns.tolist()

```

Encode categorical predictors

```

encoder = OneHotEncoder(handle_unknown='ignore', sparse_output=False)

X_train_encoded = encoder.fit_transform(X_train[categorical_feature])
X_test_encoded = encoder.transform(X_test[categorical_feature])

# Convert the encoded features back to DataFrame
X_train_encoded_df = pd.DataFrame(X_train_encoded, columns=encoder.get_feature_names_out(categorical_feature))
X_test_encoded_df = pd.DataFrame(X_test_encoded, columns=encoder.get_feature_names_out(categorical_feature))

# Concatenate the encoded features with the original numerical features
X_train_final = pd.concat([X_train_encoded_df, X_train[numerical_feature].reset_index(drop=True)], axis=1)
X_test_final = pd.concat([X_test_encoded_df, X_test[numerical_feature].reset_index(drop=True)], axis=1)

```

By default, a single decision tree grows to its full depth, which often leads to overfitting as shown below

```

# build a decision tree regressor using the default parameters
tree_reg = DecisionTreeRegressor(random_state=42)
tree_reg.fit(X_train_final, y_train)
y_pred = tree_reg.predict(X_test_final)
rmse = root_mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f"Test RMSE: {rmse:.2f}, test R^2: {r2:.2f}")

# training rmse and r2
y_train_pred = tree_reg.predict(X_train_final)
train_rmse = root_mean_squared_error(y_train, y_train_pred)
train_r2 = r2_score(y_train, y_train_pred)
print(f"Train RMSE: {train_rmse:.2f}, train R^2: {train_r2:.2f}")

# print the depth of the tree
print(f"Depth of the tree: {tree_reg.get_depth()}")
# print the number of leaves in the tree
print(f"Number of leaves in the tree: {tree_reg.get_n_leaves()}")

```

```

Test RMSE: 6219.96, test R^2: 0.87
Train RMSE: 0.00, train R^2: 1.00
Depth of the tree: 34
Number of leaves in the tree: 5925

```

As observed, the model achieves an RMSE of 0.00 and an  $R^2$  of 100% on the training data with default parameters, indicating overfitting

To address this, we've previously explored pre-pruning and post-pruning techniques. Another effective approach is bagging, which helps reduce overfitting by lowering model variance.

Next, we'll explore how bagging can improve the performance of unpruned decision trees by reducing variance

```
#Bagging the results of 10 decision trees with the default parameters to predict car price
bagging_reg = BaggingRegressor(random_state=1,
                               n_jobs=-1).fit(X_train_final, y_train)

# make predictions on the test set
y_pred_bagging = bagging_reg.predict(X_test_final)

# calculate the RMSE and R^2 score
rmse_bagging = root_mean_squared_error(y_test, y_pred_bagging)
r2_bagging = r2_score(y_test, y_pred_bagging)

print("Test RMSE with Bagging unpruned trees:", round(rmse_bagging, 2))
print("Test R^2 score with Bagging unpruned trees:", round(r2_bagging, 2))

# training RMSE and R^2 score
y_pred_train_bagging = bagging_reg.predict(X_train_final)

# calculate the RMSE and R^2 score
rmse_train_bagging = root_mean_squared_error(y_train, y_pred_train_bagging)
r2_train_bagging = r2_score(y_train, y_pred_train_bagging)

print("Train RMSE with Bagging unpruned trees:", round(rmse_train_bagging, 2))
print("Train R^2 score with Bagging unpruned trees:", round(r2_train_bagging, 2))
```

```
Test RMSE with Bagging unpruned trees: 3758.1
Test R^2 score with Bagging unpruned trees: 0.95
Train RMSE with Bagging unpruned trees: 1501.04
Train R^2 score with Bagging unpruned trees: 0.99
```

With the default settings, bagging unpruned trees improves performance, reducing the RMSE from 6219.96 to 3758.10 and increasing the  $R^2$  score from 0.87 to 0.95.

What about bagging pruned trees? Since pruning improves the performance of a single decision tree, does that mean bagging pruned trees will also outperform bagging unpruned trees? Let's find out through implementation.

Below is the [best model](#) obtained by tuning the hyperparameters of a single decision tree.

```

# fit the decision tree regressor
pruned_tree_reg = DecisionTreeRegressor(max_depth=None, min_samples_leaf=1, min_samples_split=10)
pruned_tree_reg.fit(X_train_final, y_train)

# make predictions on the test set
y_pred = pruned_tree_reg.predict(X_test_final)
# calculate the RMSE and R^2 score
rmse = root_mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

# print the RMSE and R^2 score, keep the decimal points to 2
print("Test RMSE:", round(rmse, 2))
print("Test R^2 score:", round(r2, 2))

# print the depth of the tree
print(f"Depth of the tree: {pruned_tree_reg.get_depth()}")
# print the number of leaves in the tree
print(f"Number of leaves in the tree: {pruned_tree_reg.get_n_leaves()}")

```

Test RMSE: 4726.17  
 Test R^2 score: 0.92  
 Depth of the tree: 33  
 Number of leaves in the tree: 2558

Next, let's apply bagging to these pruned trees using the default settings.

```

# Bagging the results of 10 decision trees to predict car price
bagging_reg = BaggingRegressor(estimator=pruned_tree_reg, random_state=1,
                                n_jobs=-1).fit(X_train_final, y_train)

# make predictions on the test set
y_pred_bagging = bagging_reg.predict(X_test_final)

# calculate the RMSE and R^2 score
rmse_bagging = root_mean_squared_error(y_test, y_pred_bagging)
r2_bagging = r2_score(y_test, y_pred_bagging)

print("Test RMSE with Bagging pruned trees:", round(rmse_bagging, 2))
print("Test R^2 score with Bagging pruned trees:", round(r2_bagging, 2))

```

```
# training RMSE and R^2 score
y_pred_train_bagging = bagging_reg.predict(X_train_final)

# calculate the RMSE and R^2 score
rmse_train_bagging = root_mean_squared_error(y_train, y_pred_train_bagging)
r2_train_bagging = r2_score(y_train, y_pred_train_bagging)

print("Train RMSE with Bagging pruned trees:", round(rmse_train_bagging, 2))
print("Train R^2 score with Bagging pruned trees:", round(r2_train_bagging, 2))
```

```
Test RMSE with Bagging pruned trees: 3806.7
Test R^2 score with Bagging pruned trees: 0.95
Train RMSE with Bagging pruned trees: 1868.7
Train R^2 score with Bagging pruned trees: 0.99
```

Compared to bagging the unpruned trees, the performance is slightly worse, with bagging pruned trees the RMSE is 3806.7, bagging the unpruned trees lead to RMSE 3758.1.

### Why is bagging tuned trees worse than bagging untuned trees?

In the pruned tree, limiting the maximum depth reduces variance but increases bias, as reflected by the smaller depth and fewer leaves compared to the unpruned tree. Since bagging only reduces variance and does not affect bias, applying it to pruned trees—which have slightly higher bias—results in slightly worse performance than bagging unpruned trees.

This suggests that when using bagging, we don't necessarily need to tune the hyperparameters of the base decision tree—bagging itself effectively combats overfitting by reducing variance, much like hyperparameter tuning does.

## 6.3 Bagging Doesn't Reduce Bias

Bagging high-variance models can effectively lower overall variance, as long as the individual models are not highly correlated. However, Bagging high-bias models will still produce a high-bias ensemble.

To demonstrate this, we first fit a **shallow decision tree** with `max_depth=2`, which severely underfits the data due to its limited capacity. Then, we apply **bagging** using 10 such shallow trees (default setting) as base estimators.

Since each tree has high bias, the aggregated predictions from bagging still inherit that bias. In our results, both the single shallow tree and the bagged version yield similar (and poor) performance in terms of RMSE and  $R^2$  on both the training and test sets.

This experiment shows that if your base model is too simple to capture the underlying patterns in the data, bagging will not help. To improve performance in such cases, we need to use more expressive base models or consider methods like **boosting**, which are better suited to reducing both bias and variance.

```
# Single shallow decision tree (underfitting)
shallow_tree_reg = DecisionTreeRegressor(max_depth=2, random_state=1)
shallow_tree_reg.fit(X_train_final, y_train)

# Predict and evaluate on test set
y_pred_single = shallow_tree_reg.predict(X_test_final)
rmse_single = root_mean_squared_error(y_test, y_pred_single)
r2_single = r2_score(y_test, y_pred_single)

# Predict and evaluate on training set
y_pred_train_single = shallow_tree_reg.predict(X_train_final)
rmse_train_single = root_mean_squared_error(y_train, y_pred_train_single)
r2_train_single = r2_score(y_train, y_pred_train_single)

print("Single Shallow Tree - Test RMSE:", round(rmse_single, 2))
print("Single Shallow Tree - Test R^2:", round(r2_single, 2))
print("Single Shallow Tree - Train RMSE:", round(rmse_train_single, 2))
print("Single Shallow Tree - Train R^2:", round(r2_train_single, 2))
```

```
Single Shallow Tree - Test RMSE: 11084.97
Single Shallow Tree - Test R^2: 0.58
Single Shallow Tree - Train RMSE: 10314.42
Single Shallow Tree - Train R^2: 0.6
```

Let's bag these shallow trees

```
# Bagging with 10 shallow trees
bagging_shallow = BaggingRegressor(estimator=DecisionTreeRegressor(max_depth=2, random_state=1,
                                                                    random_state=1,
                                                                    n_jobs=-1)
bagging_shallow.fit(X_train_final, y_train)

# Predict and evaluate on test set
y_pred_bagging = bagging_shallow.predict(X_test_final)
rmse_bagging = root_mean_squared_error(y_test, y_pred_bagging)
r2_bagging = r2_score(y_test, y_pred_bagging)
```



```
# Predict and evaluate on training set
y_pred_train_bagging = bagging_shallow.predict(X_train_final)
rmse_train_bagging = root_mean_squared_error(y_train, y_pred_train_bagging)
r2_train_bagging = r2_score(y_train, y_pred_train_bagging)

print("Bagged Shallow Trees - Test RMSE:", round(rmse_bagging, 2))
print("Bagged Shallow Trees - Test R^2:", round(r2_bagging, 2))
print("Bagged Shallow Trees - Train RMSE:", round(rmse_train_bagging, 2))
print("Bagged Shallow Trees - Train R^2:", round(r2_train_bagging, 2))
```

```
Bagged Shallow Trees - Test RMSE: 10894.92
Bagged Shallow Trees - Test R^2: 0.6
Bagged Shallow Trees - Train RMSE: 10114.07
Bagged Shallow Trees - Train R^2: 0.62
```

What you should observe:

- Both models show low  $R^2$  and high RMSE due to the shallow depth ( $\text{max\_depth}=2$ ).
- Bagging cannot fix the high bias inherent in a shallow decision tree.

## 6.4 Model Performance vs. Number of Trees

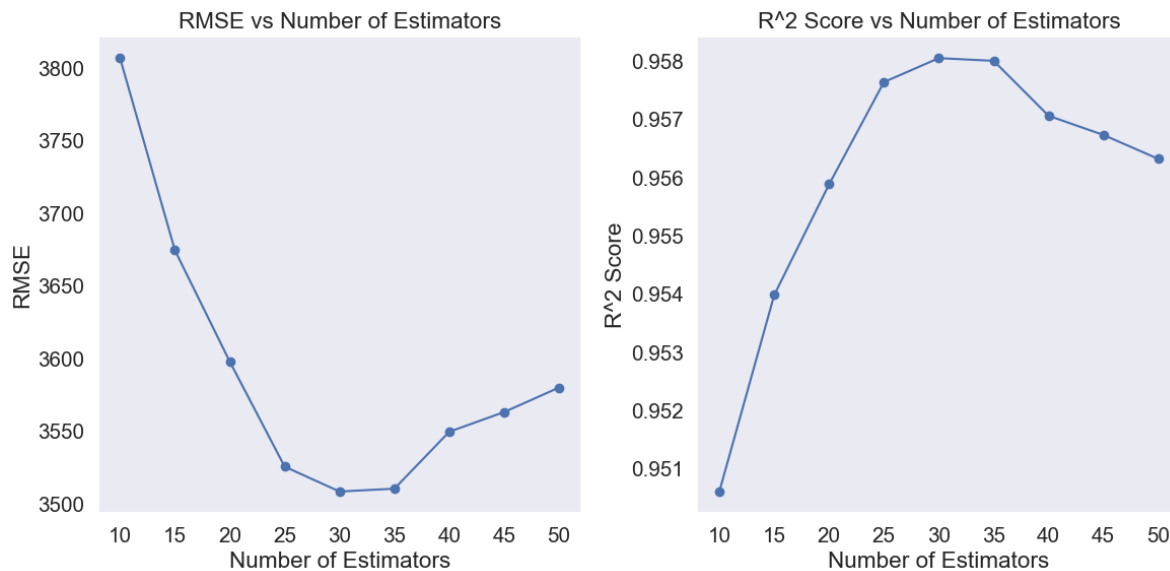
To better understand how the number of base estimators affects the performance of a bagging model, we evaluate the test RMSE and  $R^2$  score across different numbers of trees.

This analysis helps us determine whether adding more trees continues to improve performance or if the model reaches a performance plateau.

```
# explore how the number of estimators affects the performance of the model, output both oob
n_estimators = [ 10, 15, 20, 25, 30, 35, 40, 45, 50]
rmse_scores = []
r2_scores = []

# iterate through the number of estimators and fit the model
for n in n_estimators:
    bagging_reg = BaggingRegressor(estimator=pruned_tree_reg, n_estimators=n, random_state=1,
                                   n_jobs=-1).fit(X_train_final, y_train)
    y_pred_bagging = bagging_reg.predict(X_test_final)
    rmse_scores.append(np.sqrt(np.mean((y_test - y_pred_bagging) ** 2)))
    r2_scores.append(r2_score(y_test, y_pred_bagging))
```

```
# plot the RMSE and R^2 scores against the number of estimators
plt.figure(figsize=(12, 6))
plt.subplot(1, 2, 1)
plt.plot(n_estimators, rmse_scores, marker='o')
plt.title('RMSE vs Number of Estimators')
plt.xlabel('Number of Estimators')
plt.ylabel('RMSE')
plt.xticks(n_estimators)
plt.grid()
plt.subplot(1, 2, 2)
plt.plot(n_estimators, r2_scores, marker='o')
plt.title('R^2 Score vs Number of Estimators')
plt.xlabel('Number of Estimators')
plt.ylabel('R^2 Score')
plt.xticks(n_estimators)
plt.grid()
plt.tight_layout()
plt.show()
```



### Quick Takeaway

- **Increasing the number of estimators initially improves performance**, as seen from the decreasing RMSE and increasing R<sup>2</sup> scores.
- **Performance stabilizes around 30–35 estimators.** Beyond this point, additional trees provide minimal gains.

- Due to the small and possibly noisy test set, performance may appear to **peak and then decline slightly**. However, this is not typical—under normal circumstances, performance **levels off and forms a plateau**.

```
# get the number of estimators that gives the best RMSE score
best_rmse_index = np.argmin(rmse_scores)
best_rmse_n_estimators = n_estimators[best_rmse_index]
best_rmse_value = rmse_scores[best_rmse_index]
print("Best number of estimators for RMSE:", best_rmse_n_estimators)
print("Best RMSE value:", round(best_rmse_value, 2))

# get the number of estimators that gives the best R^2 score
best_r2_index = np.argmax(r2_scores)
best_r2_n_estimators = n_estimators[best_r2_index]
best_r2_value = r2_scores[best_r2_index]
print("Best number of estimators for R^2 score:", best_r2_n_estimators)
print("Best R^2 score:", round(best_r2_value, 2))
```

```
Best number of estimators for RMSE: 30
Best RMSE value: 3508.31
Best number of estimators for R^2 score: 30
Best R^2 score: 0.96
```

## 6.5 OOB Sample and OOB Score in Bagging

When training a **Bagging ensemble**, such as `BaggingClassifier` or `BaggingRegressor`, each base learner is trained on a **bootstrap sample**—a random sample *with replacement* from the original dataset.

### 6.5.1 What is an OOB Sample?

For each base learner, the data points **not selected** in the bootstrap sample form its **Out-of-Bag (OOB) sample**. On average, about **1/3 of the original data points** are not included in each bootstrap sample. These unused samples are called **OOB samples**.

### 6.5.2 What is OOB Score?

Each base learner can be evaluated on its corresponding OOB sample—i.e., the instances it did *not* see during training. This provides a **built-in validation mechanism** without needing an explicit validation set or cross-validation.

The OOB score is the **average performance** (e.g., accuracy for classifiers,  $R^2$  for regressors) of the ensemble evaluated on OOB samples.

**Note:** By default, the `oob_score` option is turned **off** in scikit-learn. You must explicitly set `oob_score=True` to enable it, as shown below.

```
from sklearn.ensemble import BaggingClassifier
```

```
bagging_clf = BaggingClassifier(
    base_estimator=DecisionTreeClassifier(),
    n_estimators=100,
    oob_score=True,
    random_state=42
)
bagging_clf.fit(X_train, y_train)
```

```
# Access the OOB score
print(f"OOB Score: {bagging_clf.oob_score_:.4f}")
```

```
n_estimators = [ 10, 15, 20, 25, 30, 35, 40, 45, 50]
rmse_scores = []
r2_scores = []
oob_scores = []
oob_rmse_scores = []
for n in n_estimators:
    bagging_reg = BaggingRegressor(estimator=tree_reg, n_estimators=n, oob_score=True, random
                                   n_jobs=-1).fit(X_train_final, y_train)
    y_pred_bagging = bagging_reg.predict(X_test_final)
    rmse_scores.append(np.sqrt(np.mean((y_test - y_pred_bagging) ** 2)))
    r2_scores.append(r2_score(y_test, y_pred_bagging))
    oob_scores.append(bagging_reg.oob_score_)
    oob_rmse_scores.append(np.sqrt(np.mean((y_train - bagging_reg.oob_prediction_) ** 2)))

# plot the RMSE and R^2 scores against the number of estimators
plt.figure(figsize=(12, 6))
plt.subplot(1, 2, 1)
plt.plot(n_estimators, rmse_scores, marker='o')
```

```

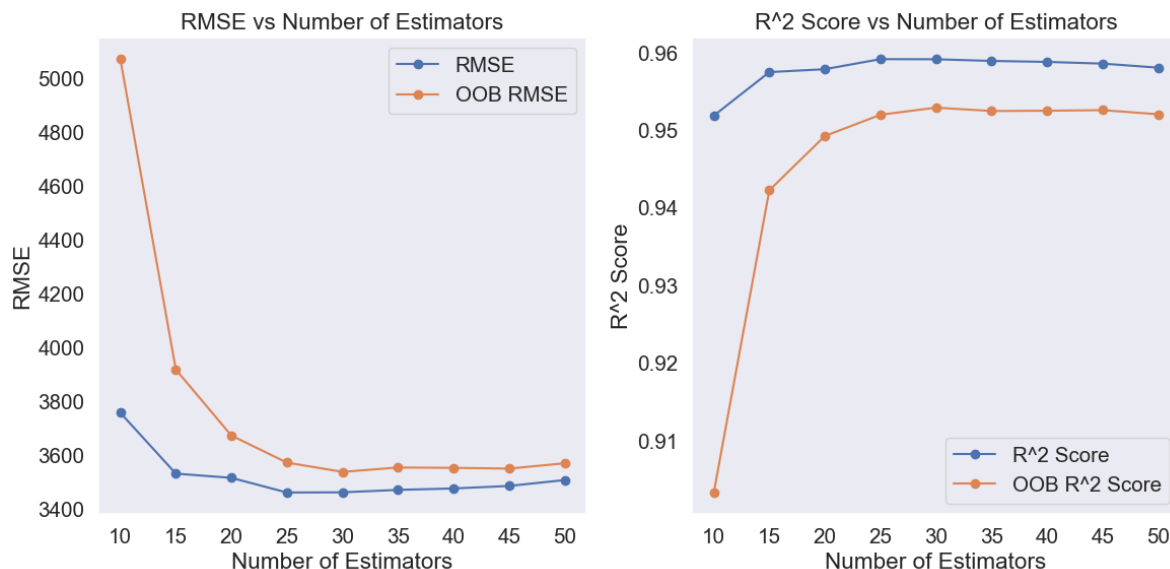
plt.plot(n_estimators, oob_rmse_scores, marker='o')
plt.legend(['RMSE', 'OOB RMSE'])
plt.title('RMSE vs Number of Estimators')
plt.xlabel('Number of Estimators')
plt.ylabel('RMSE')
plt.xticks(n_estimators)
plt.grid()
plt.subplot(1, 2, 2)
plt.plot(n_estimators, r2_scores, marker='o')
plt.plot(n_estimators, oob_scores, marker='o')
plt.legend(['R^2 Score', 'OOB R^2 Score'])
plt.title('R^2 Score vs Number of Estimators')
plt.xlabel('Number of Estimators')
plt.ylabel('R^2 Score')
plt.xticks(n_estimators)
plt.grid()
plt.tight_layout()
plt.show()

```

```

c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\sklearn\ensemble\_bagging.py:1315
warn(
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\sklearn\ensemble\_bagging.py:1315
warn(
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\sklearn\ensemble\_bagging.py:1315
warn(

```



## Quick Takeaway

- **OOB estimates become more reliable** as the number of trees grows, with OOB scores closely tracking the test performance after around 30 estimators.
- **OOB RMSE is consistently higher** and **OOB  $R^2$  is consistently lower** than their test counterparts when the ensemble is small, highlighting the **instability of OOB estimates with few trees**.

```
# get the number of estimators that gives the best OOB RMSE score
best_oob_rmse_index = np.argmin(oob_rmse_scores)
best_oob_rmse_n_estimators = n_estimators[best_oob_rmse_index]
best_oob_rmse_value = oob_rmse_scores[best_oob_rmse_index]
print("Best number of estimators for OOB RMSE:", best_oob_rmse_n_estimators)
print("Best OOB RMSE value:", round(best_oob_rmse_value, 2))
```

Best number of estimators for OOB RMSE: 30

Best OOB RMSE value: 3539.1

## 6.6 Bagging Hyperparameter Tuning

To further improve the performance of our bagging model, we can tune key hyperparameters such as:

- **n\_estimators**: the number of base estimators in the ensemble,
- **max\_features**: the maximum number of features considered at each split,
- **max\_samples**: the size of each bootstrap sample used to train base estimators,
- **bootstrap**: whether sampling is performed with replacement (**True**) or without (**False**),
- **bootstrap\_features**: whether features are sampled with replacement when selecting subsets of features for each estimator.

By systematically exploring different combinations of these parameters, we aim to identify the optimal settings that enhance predictive accuracy while maintaining good generalization.

There are two common approaches for tuning these hyperparameters: - **Cross-validation**, which provides a robust estimate of model performance, and  
- **Out-of-Bag (OOB) score**, which offers an efficient built-in alternative without needing a separate validation set.

### 6.6.1 Tuning with Cross-Validation

Next, let's use `GridSearchCV` to tune these hyperparameters and identify the best combination for improved model performance.

```
# hyperparameter tuning using GridSearchCV
param_grid = {
    'n_estimators': [10, 20, 30, 40, 50],
    'max_samples': [0.5, 0.75, 1.0],
    'max_features': [0.5, 0.75, 1.0],
    'bootstrap': [True, False],
    'bootstrap_features': [True, False]
}

bagging_reg_grid = BaggingRegressor(random_state=42, n_jobs=-1)
grid_search = GridSearchCV(bagging_reg_grid, param_grid, cv=5, scoring='neg_mean_squared_error')
grid_search.fit(X_train_final, y_train)

# get the best parameters and the best score
best_params = grid_search.best_params_
best_score = np.sqrt(-grid_search.best_score_)
print("Best parameters:", best_params)
print("Best RMSE cv score:", round(best_score, 2))

# make predictions on the test set using the best parameters
best_bagging_reg = grid_search.best_estimator_
y_pred_best_bagging = best_bagging_reg.predict(X_test_final)

# calculate the RMSE and R^2 score
rmse_best_bagging = root_mean_squared_error(y_test, y_pred_best_bagging)
r2_best_bagging = r2_score(y_test, y_pred_best_bagging)
print("Test RMSE with best Bagging model:", round(rmse_best_bagging, 2))
print("Test R^2 score with best Bagging model:", round(r2_best_bagging, 2))

# training RMSE and R^2 score
y_pred_train_best_bagging = best_bagging_reg.predict(X_train_final)

# calculate the RMSE and R^2 score
rmse_train_best_bagging = root_mean_squared_error(y_train, y_pred_train_best_bagging)
r2_train_best_bagging = r2_score(y_train, y_pred_train_best_bagging)
print("Train RMSE with best Bagging model:", round(rmse_train_best_bagging, 2))
print("Train R^2 score with best Bagging model:", round(r2_train_best_bagging, 2))
```

```
Best parameters: {'bootstrap': True, 'bootstrap_features': False, 'max_features': 0.75, 'max_
Best RMSE cv score: 3288.03
Test RMSE with best Bagging model: 3348.45
Test R^2 score with best Bagging model: 0.96
Train RMSE with best Bagging model: 1411.13
Train R^2 score with best Bagging model: 0.99
```

After simultaneously tuning multiple hyperparameters of the bagging model, including `n_estimators`, `max_features`, and `max_samples`, we achieved the best performance:

- **Test RMSE with best Bagging model: 3348.45**
- **Test R<sup>2</sup> score with best Bagging model: 0.96**

This demonstrates that careful tuning of bagging-specific parameters can lead to further improvements beyond using default or even optimized single decision trees.

## 6.6.2 Tuning with Out-of-Bag (OOB) Score

As an alternative to cross-validation, we can use the **Out-of-Bag (OOB) score** to evaluate model performance during training.

This method is more efficient for large datasets, as it avoids the need to split data or run multiple folds.

By enabling `oob_score=True`, we can monitor performance on unseen data points (those not included in each bootstrap sample) and use this score to guide hyperparameter tuning.

```
from sklearn.model_selection import ParameterGrid
# tune the hyperparameters of the decision tree regressor using oob score
# Hyperparameter grid
param_grid = {
    'n_estimators': [10, 20, 30, 40, 50],
    'max_samples': [0.5, 0.75, 1.0],
    'max_features': [0.5, 0.75, 1.0],
    'bootstrap': [True], # Required for OOB
    'bootstrap_features': [True, False]
}

# Track best parameters and OOB score
best_score = -np.inf
best_params = {}

# Iterate over all hyperparameter combinations
```



```

for params in ParameterGrid(param_grid):
    # Train model with current params and OOB score enabled
    model = BaggingRegressor(
        estimator=tree_reg,
        **params,
        oob_score=True,
        random_state=42,
        n_jobs=-1
    )
    model.fit(X_train_final, y_train)

    # Get OOB score (higher is better for R2, lower for RMSE)
    current_score = model.oob_score_

    # Update best params if current score is better
    if current_score > best_score:
        best_score = current_score
        best_params = params

# Best model
best_model = BaggingRegressor(
    estimator=tree_reg,
    **best_params,
    oob_score=True,
    random_state=42,
    n_jobs=-1
)
best_model.fit(X_train_final, y_train)

print("Best Hyperparameters:", best_params)
print("Best OOB Score:", best_score)

```

```

c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\sklearn\ensemble\_bagging.py:1315
warn(
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\sklearn\ensemble\_bagging.py:1315
warn(
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\sklearn\ensemble\_bagging.py:1315
warn(
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\sklearn\ensemble\_bagging.py:1315
warn(
c:\Users\lsi8012\AppData\Local\anaconda3\Lib\site-packages\sklearn\ensemble\_bagging.py:1315
warn(

```



### 6.6.3 Comparing Hyperparameter Tuning: Cross-Validation vs. OOB Score

Aspect	Cross-Validation (CV)	Out-of-Bag (OOB) Score
<b>Mechanism</b>	Splits training data into multiple folds to validate	Uses unused (out-of-bag) samples in each bootstrap
<b>Requires Manual Splits?</b>	Yes	No — internal to bagging process
<b>Efficiency</b>	Slower, especially for large datasets	Faster and more efficient for large datasets
<b>Bias-Variance Tradeoff</b>	More stable and less biased performance estimate	Slightly more variable and can underestimate accuracy
<b>Availability</b>	Available for all models	Only available when <code>bootstrap=True</code> in bagging
<b>Integration in Sklearn</b>	Built-in via <code>GridSearchCV</code>	Must be implemented manually for tuning
<b>Flexibility</b>	Works with any model type	Only works with bagging-based models
<b>Use Case</b>	Ideal for robust model comparison	Great for quick tuning on large datasets
<b>Scoring Access</b>	<code>.best_score_</code> from <code>GridSearchCV</code>	<code>.oob_score_</code> from trained model

#### 6.6.3.1 Best Practices for Imbalanced Classification

- **Prefer Cross-Validation**, especially with:
  - `StratifiedKFold` to maintain class distribution in each fold.
  - Custom metrics (e.g., F1-score, ROC-AUC, balanced accuracy) using `scoring=`.
- **Be cautious using OOB score**:
  - OOB score may be misleading for **rare classes**, especially when `n_estimators` is small.
  - Only use it for **rough estimates** or **early tuning** when computational efficiency is critical.

#### Summary

- **Use Cross-Validation** when:
  - You want robust, model-agnostic performance evaluation.
  - You need precise comparisons between different model types or pipelines.
  - You work on imbalanced classification task
- **Use OOB Score** when:

- You’re working with large datasets and want faster tuning.
- Your model is based on bagging (e.g., `BaggingClassifier`, `RandomForestClassifier`).
- You want to avoid manual train/validation splits.

Note: Scikit-learn’s `GridSearchCV` does **not** use **OOB score** for tuning—even if `oob_score=True` is set. To use OOB for tuning, you must loop over hyperparameters manually and evaluate using `.oob_score_`.

## 6.7 Bagging Classification Trees

Let’s revisit the same dataset used for building a single classification tree.

When using the default settings, the tree tends to **overfit** the data, as shown [here](#).

In that notebook, we addressed the overfitting issue using both **pre-pruning** and **post-pruning** techniques.

Now, we’ll explore an alternative approach—**bagging**—to reduce overfitting and improve model performance.

```
# load the dataset
heart_df = pd.read_csv('datasets/heart_disease_classification.csv')
print(heart_df .shape)
heart_df .head()
```

(303, 14)

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

```
# split the x and y data
X_cls = heart_df.drop(columns=['target'])
y_cls = heart_df.target

# split the data into train and test sets, add _cls to the variable names
X_train_cls, X_test_cls, y_train_cls, y_test_cls = train_test_split(X_cls, y_cls, test_size=0.2)
```

```
# using tree bagging to fit the data

bagging = BaggingClassifier(DecisionTreeClassifier(), n_estimators=100, random_state=0)
bagging.fit(X_train_cls, y_train_cls)

y_pred_train_cls = bagging.predict(X_train_cls)
y_pred_cls = bagging.predict(X_test_cls)

#print out the accuracy on test set and training set
print("Train Accuracy is ", accuracy_score(y_train_cls,y_pred_train_cls)*100)
print("Test Accuracy is ", accuracy_score(y_test_cls,y_pred_cls)*100)
```

```
Train Accuracy is  100.0
Test Accuracy is  85.24590163934425
```

## 7 Random Forests

*Read section 8.2.2 of the book before using these notes.*

*Note that in this course, lecture notes are not sufficient, you must read the book for better understanding. Lecture notes are just implementing the concepts of the book on a dataset, but not explaining the concepts elaborately.*

```
# Importing necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(font_scale=1.35)

# import the decision tree regressor
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier, plot_tree, export_graphviz
from sklearn.ensemble import BaggingRegressor, BaggingClassifier

# split the dataset into training and testing sets
from sklearn.model_selection import train_test_split

from sklearn.model_selection import cross_val_score, GridSearchCV, cross_val_predict, KFold

from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, FunctionTransformer
from sklearn.metrics import root_mean_squared_error, r2_score, make_scorer, accuracy_score
```

### 7.1 Motivation: Bagging Revisited

In Bagging (Bootstrap Aggregating), we:

- Train many trees on different bootstrap samples of the training data.

- Aggregate their predictions by averaging (regression) or voting (classification).

Bagging helps reduce variance. But if the trees are too similar (i.e., highly correlated), averaging won't help as much.

```
# Load the dataset
car = pd.read_csv('Datasets/car.csv')
car.head()
```

	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	vw	Beetle	2014	Manual	55457	Diesel	30	65.3266	1.6	7490
1	vauxhall	GTC	2017	Manual	15630	Petrol	145	47.2049	1.4	10998
2	merc	G Class	2012	Automatic	43000	Diesel	570	25.1172	3.0	44990
3	audi	RS5	2019	Automatic	10	Petrol	145	30.5593	2.9	51990
4	merc	X-CLASS	2018	Automatic	14000	Diesel	240	35.7168	2.3	28990

```
X = car.drop(columns=['price'])
y = car['price']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# extract the categorical columns and put them in a list
categorical_feature = X.select_dtypes(include=['object']).columns.tolist()

# extract the numerical columns and put them in a list
numerical_feature = X.select_dtypes(include=['int64', 'float64']).columns.tolist()

preprocessor = ColumnTransformer(
    transformers=[
        ('num', FunctionTransformer(), numerical_feature),
        ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_feature)
    ],
    remainder='passthrough'
)
```

### 7.1.1 Let's build a single decision tree and output its performance

```

# build a single decision tree regressor
single_tree_regressor = DecisionTreeRegressor(random_state=0)

# pipeline for the single decision tree regressor
single_tree_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('tree', single_tree_regressor)
])

# fit the pipeline to the training data
single_tree_pipeline.fit(X_train, y_train)
# make predictions on the test data
y_pred_single_tree = single_tree_pipeline.predict(X_test)
# calculate the RMSE and R^2 score
rmse_single_tree = root_mean_squared_error(y_test, y_pred_single_tree)
r2_single_tree = r2_score(y_test, y_pred_single_tree)
print(f'Single Tree RMSE: {rmse_single_tree:.2f}')
print(f'Single Tree R^2: {r2_single_tree:.2f}')

# calculate the RMSE and R^2 score for the training data
y_pred_train_single_tree = single_tree_pipeline.predict(X_train)
rmse_train_single_tree = root_mean_squared_error(y_train, y_pred_train_single_tree)
r2_train_single_tree = r2_score(y_train, y_pred_train_single_tree)
print(f'Single Tree Train RMSE: {rmse_train_single_tree:.2f}')
print(f'Single Tree Train R^2: {r2_train_single_tree:.2f}')

```

```

Single Tree RMSE: 5073.81
Single Tree R^2: 0.91
Single Tree Train RMSE: 0.00
Single Tree Train R^2: 1.00

```

```

# single tree depth
tree_depth = single_tree_pipeline.named_steps['tree'].get_depth()
print(f"Depth of the single decision tree: {tree_depth}")

```

```

Depth of the single decision tree: 34

```

### 7.1.2 Let's Build a Bagging Tree with Bootstrap Sampling to Reduce Variance

By default, `bootstrap=True`, meaning each training set is created by sampling **with replacement** from the original dataset.



```

# bagging with bootstrap
bagging_with_bootstrap_regressor = BaggingRegressor(
    estimator=DecisionTreeRegressor(random_state=0),
    n_estimators=50,
    random_state=42
)

# create a pipeline with the preprocessor and the bagging regressor
bootstrap_bagging_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('bagging', bagging_with_bootstrap_regressor)
])

# fit the pipeline to the training data
bootstrap_bagging_pipeline.fit(X_train, y_train)
# make predictions on the test data
y_pred = bootstrap_bagging_pipeline.predict(X_test)
# calculate the RMSE and R^2 score
rmse = root_mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f'RMSE using bootstrapping: {rmse:.2f}')
print(f'R^2 using bootstrapping: {r2:.2f}')
# calculate the training rmse and r^2 score
y_train_pred = bootstrap_bagging_pipeline.predict(X_train)
train_rmse = root_mean_squared_error(y_train, y_train_pred)
train_r2 = r2_score(y_train, y_train_pred)
print(f'Training RMSE using bootstrapping: {train_rmse:.2f}')
print(f'Training R^2 using bootstrapping: {train_r2:.2f}')

```

```

RMSE using bootstrapping: 3756.85
R^2 using bootstrapping: 0.95
Training RMSE using bootstrapping: 1395.75
Training R^2 using bootstrapping: 0.99

```

The test RMSE improved significantly from 5073 to 3756, and the  $R^2$  score increased from 0.91 to 0.95 after applying bagging.

### 7.1.3 Let's Build a Bagging Tree Without Bootstrap Sampling

Now, we'll turn off bootstrap sampling (`bootstrap=False`) and observe how it affects the bagging model's performance.

```

bagging_without_bootstrap_regressor = BaggingRegressor(
    estimator=DecisionTreeRegressor(random_state=0),
    bootstrap=False,
    n_estimators=50,
    random_state=42
)

# create a pipeline with the preprocessor and the bagging regressor
without_bootstrap_bagging_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('bagging', bagging_without_bootstrap_regressor)
])

# fit the pipeline to the training data
without_bootstrap_bagging_pipeline.fit(X_train, y_train)
# make predictions on the test data
y_pred = without_bootstrap_bagging_pipeline.predict(X_test)

# calculate the RMSE and R^2 score
rmse = root_mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f'RMSE without bootstrap sampling: {rmse:.2f}')
print(f'R^2 without bootstrap sampling: {r2:.2f}')

# calculate the training rmse and r^2 score
y_train_pred = without_bootstrap_bagging_pipeline.predict(X_train)
train_rmse = root_mean_squared_error(y_train, y_train_pred)
train_r2 = r2_score(y_train, y_train_pred)
print(f'Training RMSE without bootstrap sampling: {train_rmse:.2f}')
print(f'Training R^2 without bootstrap sampling: {train_r2:.2f}')

```

RMSE: 4667.43

R<sup>2</sup>: 0.93

Training RMSE: 0.00

Training R<sup>2</sup>: 1.00

As observed from the results, the performance of bagging without bootstrap sampling is worse (RMSE: 4667) compared to using bootstrap sampling (RMSE: 3756).

#### 7.1.4 Why Does Bagging Without Bootstrap Perform Worse?

When `bootstrap=True`, each tree in the ensemble is trained on a different bootstrap sample — a random sample drawn **with replacement** from the training data. This process has two key effects:

- It **introduces diversity** among the individual trees.
- It **reduces correlation** between trees.

This diversity is the **core strength** of bagging: even though individual trees may overfit, their errors tend to cancel out when their predictions are averaged, leading to improved generalization.

#### 7.1.5 Why Can Bagging Without Bootstrap Still Show Slight Improvement?

When `bootstrap=False`, all trees are trained on the **same full dataset**, removing the primary source of diversity in bagging. As a result, the **variance reduction benefit** from averaging is significantly weakened.

However, even when trees are trained on the same data, slight variations can still arise due to **internal randomness** in how decision trees are constructed. For example:

- When multiple splits yield the same information gain, one split may be selected **randomly**.
- **Ties** between split candidates can be broken differently.
- **Minor numerical differences** can occur due to floating-point operations.

These small variations cause the trees to differ slightly, allowing the ensemble to achieve **some variance reduction**, which can **slightly improve generalization** compared to a single decision tree.

However, this improvement is typically **much smaller** than the improvement achieved when using full bootstrap sampling (`bootstrap=True`).

## 7.2 Random Forest

While **diversity** is the core strength of bagging, **Random Forest** further improves upon bagging by introducing **even more diversity among the individual trees**.

The goal of Random Forest is to **further decorrelate the trees**, which leads to improved generalization and predictive performance.

### 7.2.1 Idea Behind Random Forest

Random Forest introduces an additional source of randomness:

- At each split in a tree, instead of considering **all predictors**, Random Forest randomly selects a **subset of predictors** to evaluate.

This approach:

- Increases **diversity** among the trees.
- Decreases **correlation** between trees.
- Further **reduces the variance** of the aggregated model.

**Result:** Random Forest generally achieves better performance than standard bagging, especially on high-dimensional datasets.

### 7.2.2 Key Hyperparameter Comparison

Hyperparameter	Bagging	Random Forest
<code>bootstrap</code>	Yes	Yes
<code>max_features</code>	All features considered at each split	Random subset of features at each split
<code>oob_score</code>	Often used for evaluation	Often used for evaluation
<code>n_estimators</code>	Number of trees	Number of trees

### 7.2.3 Let's Build a Random Forest Model Using the Default Settings

The `max_features` hyperparameter controls the number of features considered when searching for the best split at each node.

By default, `max_features=1.0`, meaning **all features** are considered at every split, similar to standard bagging.

```
from sklearn.ensemble import RandomForestRegressor

rf_regressor = RandomForestRegressor(
    n_estimators=50,
    random_state=42
)

# create a pipeline with the preprocessor and the bagging regressor
rf_pipeline = Pipeline(steps=[
```

```

    ('preprocessor', preprocessor),
    ('rf', rf_regressor)
])

# fit the pipeline to the training data
rf_pipeline.fit(X_train, y_train)
# make predictions on the test data
y_pred = rf_pipeline.predict(X_test)

# calculate the RMSE and R^2 score
rmse = root_mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f'RMSE using random forest: {rmse:.2f}')
print(f'R^2 using random forest: {r2:.2f}')

# calculate the training rmse and r^2 score
y_train_pred = rf_pipeline.predict(X_train)
train_rmse = root_mean_squared_error(y_train, y_train_pred)
train_r2 = r2_score(y_train, y_train_pred)
print(f'Training RMSE using random forest: {train_rmse:.2f}')
print(f'Training R^2 using random forest: {train_r2:.2f}')

```

```

RMSE: 3747.58
R^2: 0.95
Training RMSE: 1395.03
Training R^2: 0.99

```

This result is close to that of the bagging model with bootstrap sampling (RMSE: 3756 vs. 3747).

To further decorrelate the trees, we can adjust the `max_features` parameter. Reducing `max_features` limits the number of features considered at each split, which increases diversity among the trees and helps further reduce variance.

## 7.2.4 Let's Build a Random Forest Model with `sqrt` `max_features`

There are two common options to reduce the number of features considered at each split: `sqrt` and `log2`.

Here, we will set `max_features='sqrt'` and observe how it affects the model's performance.

```

rf_sqrt_regressor = RandomForestRegressor(
    n_estimators=50,
    max_features='sqrt',
    random_state=42
)

# create a pipeline with the preprocessor and the bagging regressor
rf_sqrt_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('bagging', rf_sqrt_regressor)
])

# fit the pipeline to the training data
rf_sqrt_pipeline.fit(X_train, y_train)
# make predictions on the test data
y_pred = rf_sqrt_pipeline.predict(X_test)

# calculate the RMSE and R^2 score
rmse = root_mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f'RMSE: {rmse:.2f}')
print(f'R^2: {r2:.2f}')

# calculate the training rmse and r^2 score
y_train_pred = rf_sqrt_pipeline.predict(X_train)
train_rmse = root_mean_squared_error(y_train, y_train_pred)
train_r2 = r2_score(y_train, y_train_pred)
print(f'Training RMSE: {train_rmse:.2f}')
print(f'Training R^2: {train_r2:.2f}')

```

RMSE: 3424.28

R<sup>2</sup>: 0.96

Training RMSE: 1279.85

Training R<sup>2</sup>: 0.99

By using `sqrt` for `max_features`, we further decorrelate the trees, resulting in a lower RMSE of 3424 compared to 3747 when using all features at each split.

## 7.3 Let's Explore How max\_features Affects Performance

The `max_features` parameter controls the degree of feature decorrelation among trees. Let's explore different values:

```
# explore how the max_features parameter affects the model performance
max_features = ['sqrt', 'log2', 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0]
rmse_list = []
r2_list = []
for max_feature in max_features:
    rf_regressor = RandomForestRegressor(
        n_estimators=50,
        max_features=max_feature,
        random_state=42
    )

    # create a pipeline with the preprocessor and the bagging regressor
    rf_pipeline = Pipeline(steps=[
        ('preprocessor', preprocessor),
        ('bagging', rf_regressor)
    ])

    # fit the pipeline to the training data
    rf_pipeline.fit(X_train, y_train)
    # make predictions on the test data
    y_pred = rf_pipeline.predict(X_test)

    # calculate the RMSE and R^2 score
    rmse = root_mean_squared_error(y_test, y_pred)
    r2 = r2_score(y_test, y_pred)
    rmse_list.append(rmse)
    r2_list.append(r2)

# plot the RMSE and R^2 score against the max_features parameter
plt.figure(figsize=(12, 6))
plt.subplot(1, 2, 1)
plt.plot(max_features, rmse_list, marker='o')
plt.xlabel('max_features')
plt.ylabel('RMSE')
plt.title('RMSE vs max_features')
plt.grid(True)
plt.subplot(1, 2, 2)
```

```
plt.plot(max_features, r2_list, marker='o')
plt.xlabel('max_features')
plt.ylabel('R^2')
plt.title('R^2 vs max_features')
plt.grid(True)
plt.tight_layout()
plt.show()
```



As observed from the plots, **Random Forest performs best when tree decorrelation is balanced.**

Setting `max_features` too low or too high hurts the model's generalization ability.

- **$R^2$**  peaks when `max_features` is around **0.5 to 0.6**, consistent with the lowest RMSE values.
- **$R^2$**  drops at both extremes:
  - Using **too few features** (`sqrt`, `log2`, or very small proportions) leads to **underfitting**.
  - Using **all features** (`max_features=1.0`) increases correlation between trees, leading to **overfitting**.

**Key takeaway:** Carefully tuning `max_features` is critical for achieving the best balance between bias and variance in Random Forest models.

Let's get the minimum RMSE and the corresponding `max_features` value from the result



```

min_rmse = min(rmse_list)
min_rmse_index = rmse_list.index(min_rmse)
best_max_feature = max_features[min_rmse_index]
print(f'Minimum RMSE: {min_rmse:.2f}')
print(f'Best max_features: {best_max_feature}')
# get the maximum R^2 and the corresponding max_features parameter
max_r2 = max(r2_list)
max_r2_index = r2_list.index(max_r2)
best_max_feature_r2 = max_features[max_r2_index]
print(f'Maximum R^2: {max_r2:.2f}')
print(f'Best max_features: {best_max_feature_r2}')

```

```

Minimum RMSE: 3338.02
Best max_features: 0.5
Maximum R^2: 0.96
Best max_features: 0.5

```

Adjusting `max_features` from 'sqrt' to 0.5 led to a slight improvement in RMSE, reducing it from 3424 to 3338

## 7.4 Other Hyperparameters in Random Forest

### 7.4.1 Why Bagging Uses Unpruned Trees

- Bagging's main strength lies in **reducing variance**, not bias.
- Deep, unpruned decision trees tend to **overfit** (high variance), but bagging effectively reduces this variance through aggregation.
- Using pruned trees reduces variance but **increases bias** — and since bagging does not correct bias, this would weaken overall performance.

Therefore, in bagging, it is common to let each tree **grow fully** to preserve low bias and rely on bagging to reduce variance.

### 7.4.2 Hyperparameters That Control Tree Complexity in Random Forest

Setting	Effect	Applies to
<code>max_depth=None</code>	Full trees, low bias, high variance	Bagging, RF

Setting	Effect	Applies to
<code>max_depth=some int</code>	Pruned trees, more bias, less variance	Especially helpful in RF
<code>min_samples_split/leaves</code>	Prevents small, unreliable branches	Both

### 7.4.3 Why Random Forest Often Limits Tree Depth

In Random Forest, only a **subset of features** is considered at each split.

As a result, fully growing trees without any depth constraint can cause them to **overfit** to noise within these smaller subsets.

Limiting tree complexity in Random Forest:

- **Prevents deep trees** from chasing noise and overfitting.
- **Improves generalization**, especially in high-dimensional or noisy datasets.
- **Balances** the bias-variance tradeoff more effectively than using full trees.

Careful tuning of tree depth and other complexity-controlling hyperparameters is critical to maximizing Random Forest performance.

### 7.4.4 Let's Tune Multiple Hyperparameters Simultaneously Using Cross-Validation

Given the number of hyperparameters involved, we will use `BayesSearchCV` to efficiently perform tuning.

This approach helps reduce computational cost while exploring a wide range of hyperparameter combinations.

```
# hyperparameter tuning for the random forest regressor

from skopt.space import Integer, Categorical, Real
from skopt import BayesSearchCV

# Rename the pipeline step for clarity (recommended)
random_forest_regressor = RandomForestRegressor(
    random_state=42
)

random_forest_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('rf', random_forest_regressor) # Renamed from 'bagging' to 'rf'
```

```

])

param_space = {
    # Tree structure (control complexity)
    "rf__max_depth": Integer(5, 35),
    "rf__min_samples_split": Integer(2, 20),
    "rf__min_samples_leaf": Integer(1, 10),
    "rf__max_features": Real(0.1, 1.0),

    # Ensemble settings
    "rf__n_estimators": Integer(20, 60),

    # Advanced
    "rf__max_samples": Real(0.1, 1.0),
}

opt = BayesSearchCV(
    randome_forest_pipeline,
    param_space,
    n_iter=50, # Adjust based on computational resources
    cv=5,
    n_jobs=-1,
    random_state=42
)

opt.fit(X_train, y_train)

# make predictions on the test data
y_pred = opt.predict(X_test)
# calculate the RMSE and R^2 score
rmse = root_mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f'RMSE: {rmse:.2f}')
print(f'R^2: {r2:.2f}')
# calculate the training rmse and r^2 score
y_train_pred = opt.predict(X_train)
train_rmse = root_mean_squared_error(y_train, y_train_pred)
train_r2 = r2_score(y_train, y_train_pred)
print(f'Training RMSE: {train_rmse:.2f}')
print(f'Training R^2: {train_r2:.2f}')

```

RMSE: 3278.42

R^2: 0.96

Training RMSE: 1220.88

Training  $R^2$ : 0.99

As observed in the results, RMSE was further reduced after simultaneously tuning multiple hyperparameters compared to only tuning `max_features` (from 3338 to 3278).

However, due to the small and simple nature of the dataset, the performance improvement is relatively marginal.

On larger and more complex datasets, the performance gains from comprehensive hyperparameter tuning would likely be much more substantial.

## 7.5 Feature Importance in Random Forest

Random Forest provides a natural way to estimate **feature importance**.

Each time a feature is used to split a node, it contributes to reducing impurity (such as Gini impurity for classification or variance for regression).

By averaging these contributions over all trees in the forest, we can rank features by how important they are to the model's predictive performance.

### 7.5.1 How Feature Importance Is Calculated

- A feature's importance is based on the **total reduction of the criterion** (e.g., variance for regression) it brings across all splits it is used in.
- Features that result in larger reductions in impurity are assigned **higher importance scores**.
- The importance scores are **normalized** so that they sum to 1 across all features.

### 7.5.2 Accessing Feature Importances

After fitting a Random Forest model, feature importances can be accessed through the attribute:

```
model.feature_importances_
```

```
# get numerical_feature and categorical_feature from the pipeline
num_features = numerical_feature
cat_transformer = opt.best_estimator_.named_steps['preprocessor'].named_transformers_['cat']
cat_features = cat_transformer.get_feature_names_out(categorical_feature)
```

```

# concatenate all feature names
feature_names = np.concatenate([num_features, cat_features])

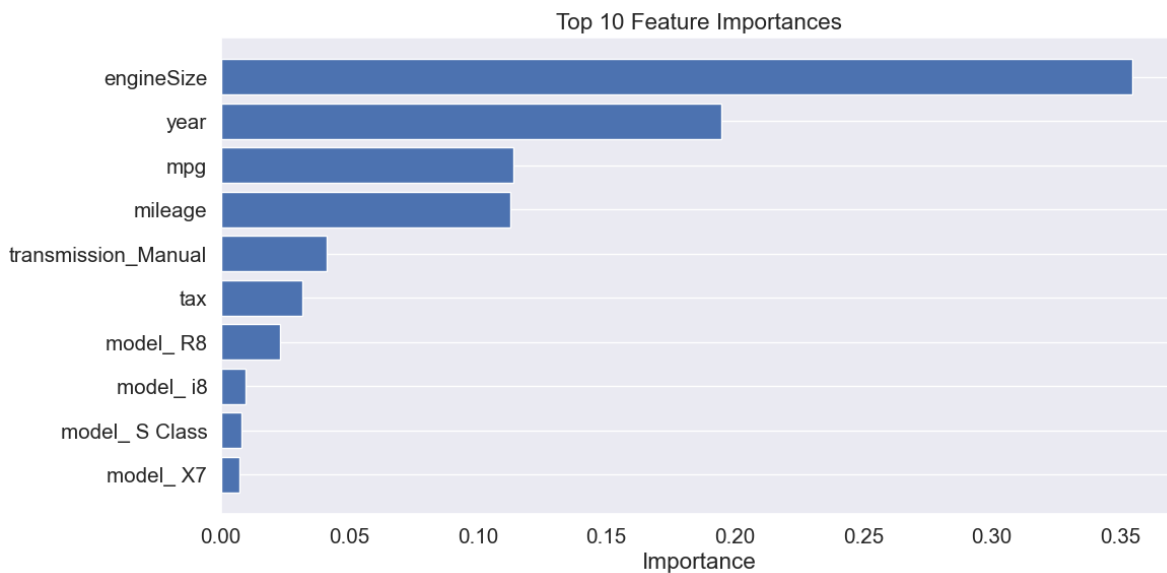
# output feature importances

importances = opt.best_estimator_.named_steps['rf'].feature_importances_
feature_importances = pd.DataFrame(importances, index=feature_names, columns=['importance'])

# select top 10 features
top_10 = feature_importances.head(10)

# plot the top 10 feature importances
plt.figure(figsize=(12, 6))
plt.barh(top_10.index[::-1], top_10['importance'][:,::-1]) # reverse for top-to-bottom order
plt.xlabel('Importance')
plt.title('Top 10 Feature Importances')
plt.grid(axis='x')
plt.tight_layout()
plt.show()

```



## 7.6 In Summary

Random Forest is a special case of bagging.

The `n_estimators` and `oob_score` hyperparameters function similarly in both methods, help-

ing to aggregate multiple decision trees into a strong ensemble.

In this notebook, we focused on the key differences between **Random Forest** and **standard bagging**.

Random Forest generally outperforms bagging by introducing an additional layer of randomness:

at each split, only a **random subset of features** is considered.

This strategy **decorrelates** the individual trees, **increases diversity** within the ensemble, and **further reduces variance**, leading to stronger generalization performance.

**Key takeaway:** Bagging reduces variance by aggregating independent models, while Random Forest improves further by strategically injecting feature-level randomness to strengthen ensemble diversity.

## 7.7 Next Step

In the next section, we will explore **Boosting methods**, where models are built **sequentially**, each one focusing on correcting the errors made by the previous models.

Boosting shifts the focus from reducing variance to **reducing bias**, offering another powerful strategy for improving model performance.

## 8 Adaptive Boosting

After learning how Bagging and Random Forest reduce variance by aggregating many trees, we now turn to a different strategy: **Boosting**.

**Boosting** builds trees **sequentially**, with each new tree focusing on correcting the mistakes made by the previous ones.

*Read section 8.2.3 of the book before using these notes.*

For the exact algorithms underlying the AdaBoost algorithm, check out the papers [AdaBoostRegressor\(\)](#) and [AdaBoostClassifier\(\)](#).

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import cross_val_score, train_test_split, KFold, cross_val_predict
from sklearn.metrics import root_mean_squared_error, r2_score, roc_curve, auc, precision_recall_curve,
recall_score, precision_score, confusion_matrix
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV, ParameterGrid, StratifiedKFold
from sklearn.ensemble import BaggingRegressor, BaggingClassifier, AdaBoostRegressor, AdaBoostClassifier,
RandomForestRegressor
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, FunctionTransformer
import itertools as it
import time as time

from skopt import BayesSearchCV
from skopt.space import Real, Categorical, Integer
from skopt.plots import plot_objective, plot_histogram, plot_convergence
import warnings
from IPython import display
```

## 8.1 What is AdaBoost?

**AdaBoost** stands for **Adaptive Boosting**.

It was one of the first boosting algorithms developed.

The core idea behind AdaBoost:

- Train a **weak learner** (usually a shallow decision tree) on the original data.
- Increase the weights of examples that the learner misclassified.
- Train the next learner on this updated, reweighted data.
- Repeat this process, focusing more and more on hard-to-predict examples.

The final prediction is a **weighted combination** of all the weak learners.

## 8.2 AdaBoost Intuition

- Easy-to-classify points are **de-emphasized**.
- Hard-to-classify points are **emphasized**.
- Each learner **adapts** to the mistakes made by previous learners — hence “adaptive” boosting.
- **Better-performing learners** are given **higher weight** in the final prediction.

Over time, the model becomes better at handling difficult cases.

## 8.3 How AdaBoost Works (High-Level Steps)

1. Initialize equal weights for all training examples.
2. Train a weak learner (e.g., decision stump).
3. Evaluate its performance:
  - Increase weights for misclassified points.
  - Decrease weights for correctly classified points.
4. Train the next learner using the updated weights.
5. Repeat for a set number of learners (**n\_estimators**).
6. Combine all learners into a final weighted model.

## 8.4 Key Hyperparameters in AdaBoost



Hyperparameter	Meaning	Typical Values
<code>n_estimators</code>	Number of weak learners	50–500
<code>learning_rate</code>	Shrinks each learner's contribution	0.01–1.0
<code>estimator</code>	Type of weak learner (default: decision stump)	Shallow trees

- **Lowering `learning_rate`** typically requires **more estimators** but improves generalization.

## 8.5 AdaBoost for Regression

We will revisit the car dataset we used earlier and evaluate how AdaBoost performs compared to a single decision tree and a bagging ensemble

```
# Load the dataset
car = pd.read_csv('Datasets/car.csv')
car.head()
```

	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	vw	Beetle	2014	Manual	55457	Diesel	30	65.3266	1.6	7490
1	vauxhall	GTC	2017	Manual	15630	Petrol	145	47.2049	1.4	10998
2	merc	G Class	2012	Automatic	43000	Diesel	570	25.1172	3.0	44990
3	audi	RS5	2019	Automatic	10	Petrol	145	30.5593	2.9	51990
4	merc	X-CLASS	2018	Automatic	14000	Diesel	240	35.7168	2.3	28990

```
print(car.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7632 entries, 0 to 7631
Data columns (total 10 columns):
#   Column          Non-Null Count  Dtype
---  -
0   brand           7632 non-null   object
1   model           7632 non-null   object
2   year            7632 non-null   int64
3   transmission     7632 non-null   object
4   mileage          7632 non-null   int64
5   fuelType        7632 non-null   object
6   tax              7632 non-null   int64
```

```

7   mpg          7632 non-null   float64
8   engineSize   7632 non-null   float64
9   price        7632 non-null   int64
dtypes: float64(2), int64(4), object(4)
memory usage: 596.4+ KB
None

```

```

X = car.drop(columns=['price'])
y = car['price']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# extract the categorical columns and put them in a list
categorical_feature = X.select_dtypes(include=['object']).columns.tolist()

# extract the numerical columns and put them in a list
numerical_feature = X.select_dtypes(include=['int64', 'float64']).columns.tolist()

```

```

encoder = OneHotEncoder(handle_unknown='ignore', sparse_output=False)

X_train_encoded = encoder.fit_transform(X_train[categorical_feature])
X_test_encoded = encoder.transform(X_test[categorical_feature])

# Convert the encoded features back to DataFrame
X_train_encoded_df = pd.DataFrame(X_train_encoded, columns=encoder.get_feature_names_out(categorical_feature))
X_test_encoded_df = pd.DataFrame(X_test_encoded, columns=encoder.get_feature_names_out(categorical_feature))

# Concatenate the encoded features with the original numerical features
X_train_final = pd.concat([X_train_encoded_df, X_train[numerical_feature]].reset_index(drop=True))
X_test_final = pd.concat([X_test_encoded_df, X_test[numerical_feature]].reset_index(drop=True))

```

### 8.5.1 Let's build a adaboost regressor model with default settings

```

# build a adaboost regressor model with default parameters
adaboost_regressor = AdaBoostRegressor(random_state=0)

# fit the model
adaboost_regressor.fit(X_train_final, y_train)

```

```

# predict the test set
y_pred = adaboost_regressor.predict(X_test_final)

# calculate the mean squared error
rmse = root_mean_squared_error(y_test, y_pred)
print(f'RMSE: {rmse:.2f}')
print(f'R2 Score: {r2_score(y_test, y_pred)}')

# calculate the RMSE and R^2 score for the training data
y_pred_train = adaboost_regressor.predict(X_train_final)
rmse_train = root_mean_squared_error(y_train, y_pred_train)
r2_train = r2_score(y_train, y_pred_train)
print(f'Adaboost Train RMSE: {rmse_train:.2f}')
print(f'Adaboost Train R^2: {r2_train:.2f}')

# calculate the test score

```

```

RMSE: 10239.74
R2 Score: 0.6426025126917443
Adaboost Train RMSE: 10081.26
Adaboost Train R^2: 0.62

```

#### 8.5.1.1 Why AdaBoost Perform Worse Here

- **Default base estimator is very weak:**  
By default, AdaBoost uses **Decision Stumps** (`DecisionTreeRegressor(max_depth=1)`), which are extremely shallow and tend to **underfit** the data badly.
- **Learning rate (`learning_rate=1.0`) is too aggressive:**  
When using very weak learners, a high learning rate can cause the boosting process to **fail to properly build up model strength**, leading to poor performance.
- **Dataset characteristics:**  
This dataset is **small** and **not very noisy**, using very shallow trees combined with a high learning rate can cause **severe underfitting**.

What should we do next to reduce bias

- Use deeper Trees as Base Learners
- Tune Learning Rate
- Increase the Number of Estimators

### 8.5.2 Impact of Tree Depth

By default, AdaBoost uses shallow decision stumps (`max_depth=1`) as weak learners. However, slightly increasing the tree depth can make each learner more expressive, helping the ensemble capture more complex patterns in the data.

This often leads to a **reduction in cross-validation RMSE** and improved model performance, especially when the underlying relationships in the data are non-linear.

From our previous exploration, we found that the fully grown decision tree has a depth of 34 on this dataset.

In this section, we'll experiment with limiting the tree depth and observe how it affects the model's RMSE.

The goal is to find a depth that balances **bias and variance**, leading to better generalization.

```
# get a list of models to evaluate
def get_models():
    models = dict()
    # explore depths from 1 to 10
    for i in range(1,34):
        # define base model
        base = DecisionTreeRegressor(max_depth=i)
        # define ensemble model
        models[str(i)] = AdaBoostRegressor(estimator=base,n_estimators=100, learning_rate=0.1)
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=10, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = -cross_val_score(model, X, y, scoring='neg_root_mean_squared_error', cv=cv, n_jobs=4)
    return scores

# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X_train_final, y_train)
```

```

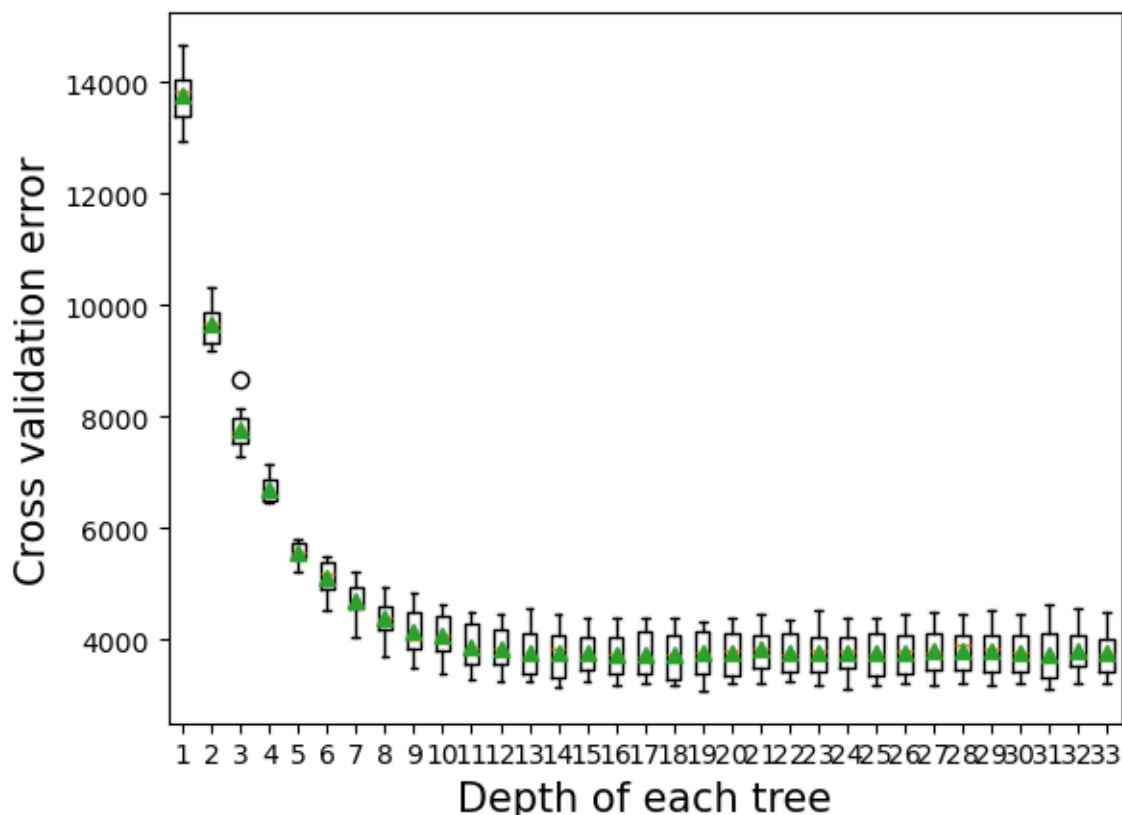
# store the results
results.append(scores)
names.append(name)
# summarize the performance along the way
print('>%s %.3f (%.3f)' % (name, np.mean(scores), np.std(scores)))
plt.boxplot(results, labels=names, showmeans=True)
plt.ylabel('Cross validation error',fontsize=15)
plt.xlabel('Depth of each tree',fontsize=15);

```

```

>1 13770.872 (518.361)
>2 9673.586 (398.116)
>3 7783.875 (393.200)
>4 6686.293 (253.234)
>5 5575.918 (176.859)
>6 5106.235 (342.995)
>7 4695.491 (353.584)
>8 4395.372 (366.340)
>9 4143.296 (422.938)
>10 4064.871 (414.840)
>11 3880.994 (433.105)
>12 3831.714 (396.971)
>13 3773.891 (437.190)
>14 3771.442 (425.043)
>15 3769.082 (388.875)
>16 3741.356 (394.563)
>17 3740.153 (415.062)
>18 3721.954 (444.760)
>19 3765.976 (425.594)
>20 3777.496 (425.025)
>21 3827.491 (434.580)
>22 3761.119 (409.760)
>23 3773.776 (429.259)
>24 3763.150 (408.907)
>25 3763.396 (417.512)
>26 3782.507 (405.136)
>27 3791.885 (446.067)
>28 3812.705 (406.567)
>29 3792.121 (445.264)
>30 3780.123 (429.651)
>31 3739.596 (472.987)
>32 3797.480 (426.929)
>33 3754.727 (417.368)

```



As shown in the plot, very shallow trees (e.g., `max_depth=1` to 3) result in high cross-validation error due to underfitting.

As tree depth increases, the model becomes more expressive, and the error drops sharply up to around depth 10.

Beyond this point, deeper trees offer diminishing returns, and performance stabilizes.

> This suggests that **slightly deeper trees (e.g., depth 5–10)** strike a good balance between model complexity and generalization in AdaBoost.

### 8.5.3 Impact of Learning Rate

In boosting algorithms such as **AdaBoost** or **Gradient Boosting**, the `learning_rate` controls **how much each new tree contributes** to the overall model.

Each new tree makes a correction to the current prediction, and the learning rate scales **how aggressively** that correction is applied.

- A **high learning rate** takes **large correction steps** — fast learning, but higher risk of overshooting or overfitting.

- A **low learning rate** takes **small correction steps** — more stable, but may underfit unless paired with enough trees.

### 8.5.3.1 Effect on Performance

Learning Rate	Behavior	Risk
Very Small (e.g., 0.01)	Learns slowly, needs many trees	Underfitting if not enough trees
Moderate (e.g., 0.1–0.2)	Balanced correction, stable learning	Often optimal
Large (e.g., 0.5–1.0)	Learns quickly, may overshoot	Overfitting or unstable learning

**Key takeaway:** Small learning rates usually generalize better — especially when combined with more estimators.

```
def get_models():
    models = dict()
    learning_rates = [0.01, 0.02, 0.04, 0.08, 0.1, 0.15, 0.2, 0.3, 0.6, 1.0]
    for i in range(len(learning_rates)):
        key = learning_rates[i]
        models[key] = AdaBoostRegressor(learning_rate=learning_rates[i])
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=10, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = -cross_val_score(model, X, y, scoring='neg_root_mean_squared_error', cv=cv, n_jobs=4)
    return scores

# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X_train_final, y_train)
    # store the results
    results.append(scores)
```

```

    names.append(name)
    # summarize the performance along the way
    print('>%s %.1f (%.1f)' % (name, np.mean(scores), np.std(scores)))
# plot model performance for comparison
plt.figure(figsize=(7, 7))
plt.boxplot(results, labels=names, showmeans=True)
plt.ylabel('Cross validation error',fontsize=15)
plt.xlabel('Learning rate',fontsize=15);

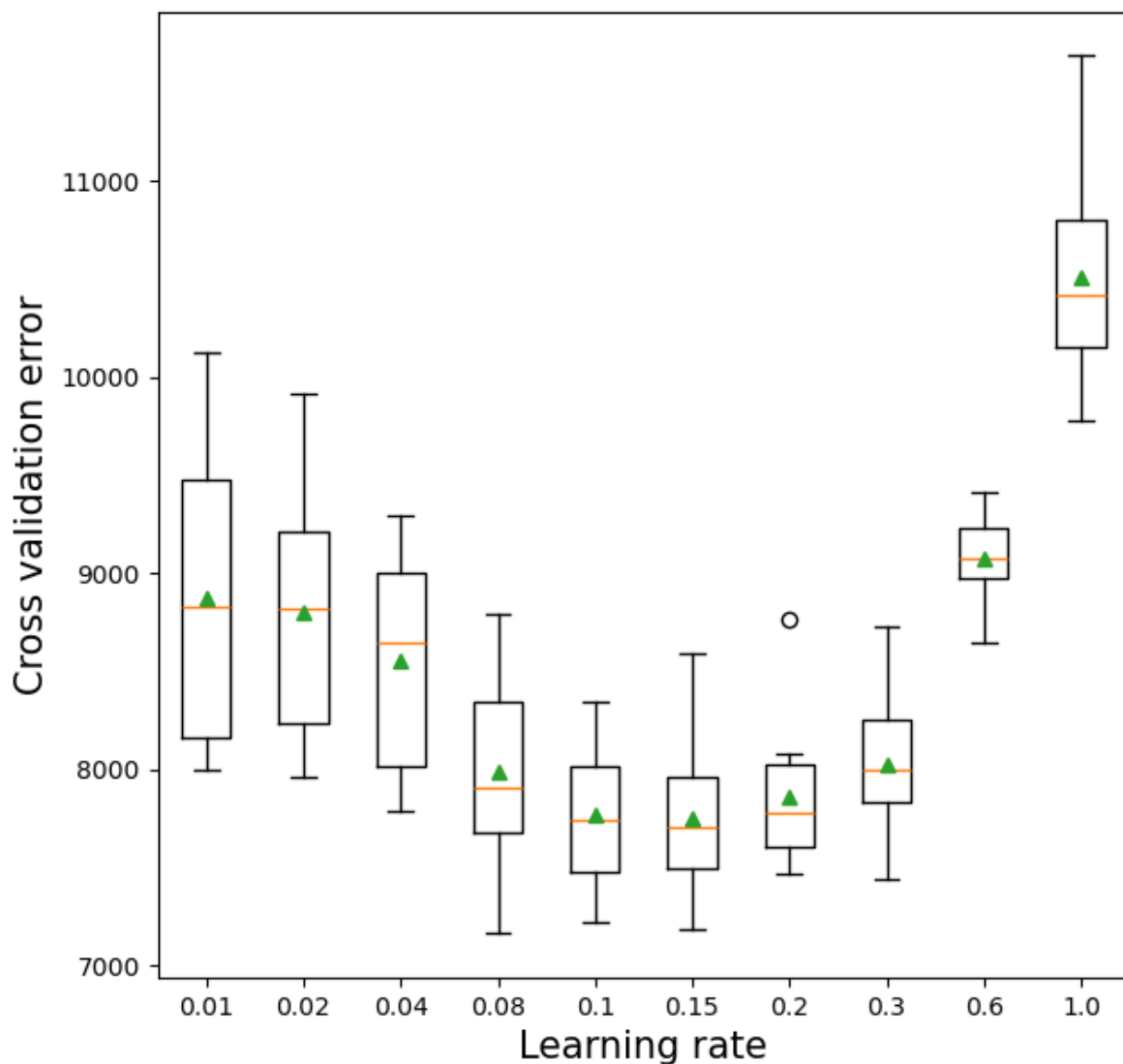
```

```

>0.01 8877.6 (725.1)
>0.02 8797.2 (656.7)
>0.04 8554.4 (540.3)
>0.08 7988.1 (479.9)
>0.1 7763.7 (345.8)
>0.15 7754.0 (380.3)
>0.2 7862.9 (368.8)
>0.3 8024.6 (345.3)
>0.6 9078.9 (205.6)
>1.0 10508.4 (507.6)

```





The plot shows that moderate learning rates (0.1–0.2) yield the best and most stable model performance, while very small or very large values hurt generalization — likely due to underfitting or overfitting.

#### 8.5.4 Impact of Number of Trees in Boosting

As the number of trees increases in a boosting model, the **prediction bias tends to decrease**, while the **variance may increase**.

This creates a trade-off:

- Too few trees → underfitting (high bias)
- Too many trees → potential overfitting (high variance)

There is typically an **optimal number of trees** that minimizes the overall prediction error, which can be identified using cross-validation.

```
def get_models():
    models = dict()
    # define number of trees to consider
    n_trees = [10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 200, 300, 500]
    for n in n_trees:
        models[str(n)] = AdaBoostRegressor(n_estimators=n, random_state=1, learning_rate=0.1)
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=5, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = -cross_val_score(model, X, y, scoring='neg_root_mean_squared_error', cv=cv, n_jobs=-1)
    return scores

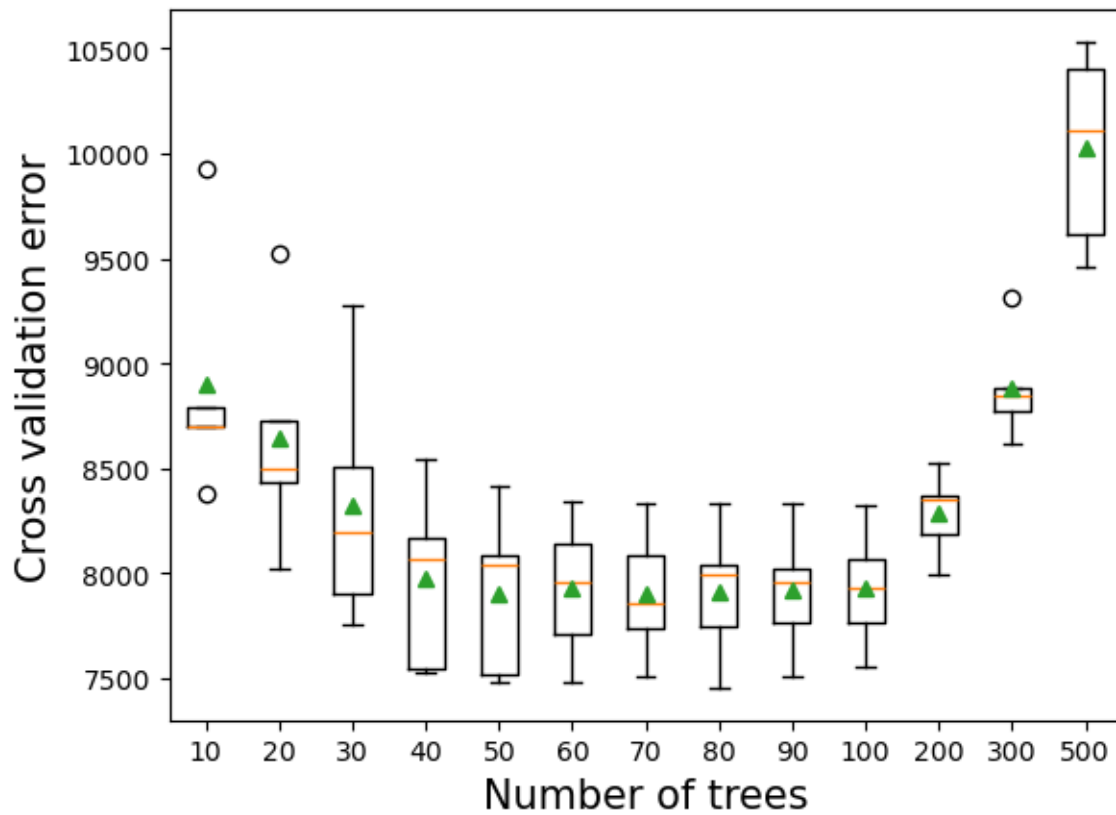
# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X_train_final, y_train)
    # store the results
    results.append(scores)
    names.append(name)
    # summarize the performance along the way
    print('>%s %.3f (%.3f)' % (name, np.mean(scores), np.std(scores)))
# plot model performance for comparison
plt.boxplot(results, labels=names, showmeans=True)
plt.ylabel('Cross validation error', fontsize=15)
plt.xlabel('Number of trees', fontsize=15);
```

```
>10 8901.126 (529.620)
>20 8640.382 (495.164)
```

```

>30 8328.349 (539.563)
>40 7972.809 (387.803)
>50 7907.280 (359.779)
>60 7927.212 (305.995)
>70 7904.131 (281.108)
>80 7914.196 (295.777)
>90 7917.841 (274.357)
>100 7927.393 (260.542)
>200 8286.386 (180.913)
>300 8884.444 (230.006)
>500 10024.047 (421.340)

```



With a learning rate of 0.1, the validation error initially decreases, then levels off, and eventually starts to increase — indicating that overfitting is beginning to occur

### 8.5.5 Tuning Hyperparameters Simultaneously

In the following section, we will use **BayesSearchCV** instead of **GridSearchCV** to efficiently tune multiple hyperparameters at once.

Unlike grid search, which exhaustively evaluates all combinations, Bayesian optimization intelligently explores the hyperparameter space by learning from previous evaluations.

This allows us to find a high-performing model using **fewer iterations and less computation**.

```
from skopt import BayesSearchCV
from skopt.space import Real, Integer

# Define the base estimator search space (DecisionTreeRegressor)
base_estimator = DecisionTreeRegressor()

# AdaBoost model (wrapped for BayesSearchCV)
adaboost = AdaBoostRegressor(estimator=base_estimator, random_state=42)

# Search space for tuning
search_space = {
    'estimator__max_depth': Integer(5, 25),
    'n_estimators': Integer(100, 500),
    'learning_rate': Real(0.01, 2.0, prior='log-uniform')
}

# Set up the BayesSearchCV
opt = BayesSearchCV(
    estimator=adaboost,
    search_spaces=search_space,
    n_iter=50,
    scoring='neg_root_mean_squared_error', # or use 'r2'
    cv=10,
    random_state=42,
    n_jobs=-1,
    verbose=1
)

# Fit on training data
opt.fit(X_train_final, y_train)

# Best parameters
print("Best parameters found:")
print(opt.best_params_)
```



```

Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Fitting 10 folds for each of 1 candidates, totalling 10 fits
Best parameters found:
OrderedDict({'estimator__max_depth': 16, 'learning_rate': 1.3460276374020355, 'n_estimators':

```

```

# Best score
print("Best score (RMSE):")
print(-opt.best_score_)

```

```

Best score (RMSE):
3280.273309604182

```

```

# evaluate the best model on the test set
best_model = opt.best_estimator_
y_pred_test = best_model.predict(X_test_final)
rmse_test = root_mean_squared_error(y_test, y_pred_test)
print(f'Test RMSE: {rmse_test:.2f}')
print(f'Test R^2: {r2_score(y_test, y_pred_test):.2f}')

```

```

Test RMSE: 3989.52
Test R^2: 0.95

```

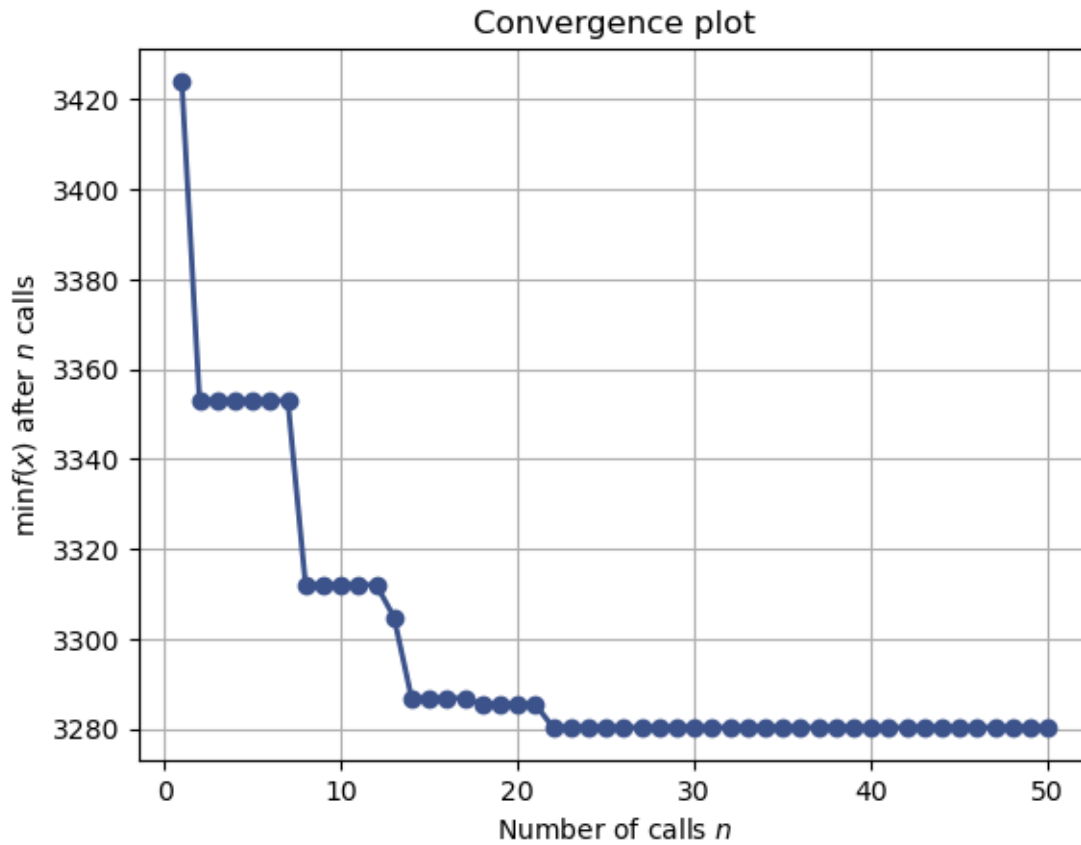
Below is the plot showing the minimum cross-validated score computed obtained until 'n' hyperparameter values are considered for cross-validation.

```

# import plot_convergence from skopt
from skopt.plots import plot_convergence

plot_convergence(opt.optimizer_results_)
plt.show()

```



```
# access the full results
results_df = pd.DataFrame(opt.cv_results_)
results_df['mean_test_score'] = -results_df['mean_test_score']
results_df.head()
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_estimator__max_depth
0	27.288841	0.743828	0.243362	0.085644	13
1	17.667935	0.272929	0.097709	0.037637	22
2	8.346380	0.144302	0.054507	0.011088	14
3	30.282255	0.495589	0.178685	0.054628	21
4	27.381998	0.592005	0.195445	0.069875	21

### 8.5.5.1 Analyzing BayesSearchCV Results

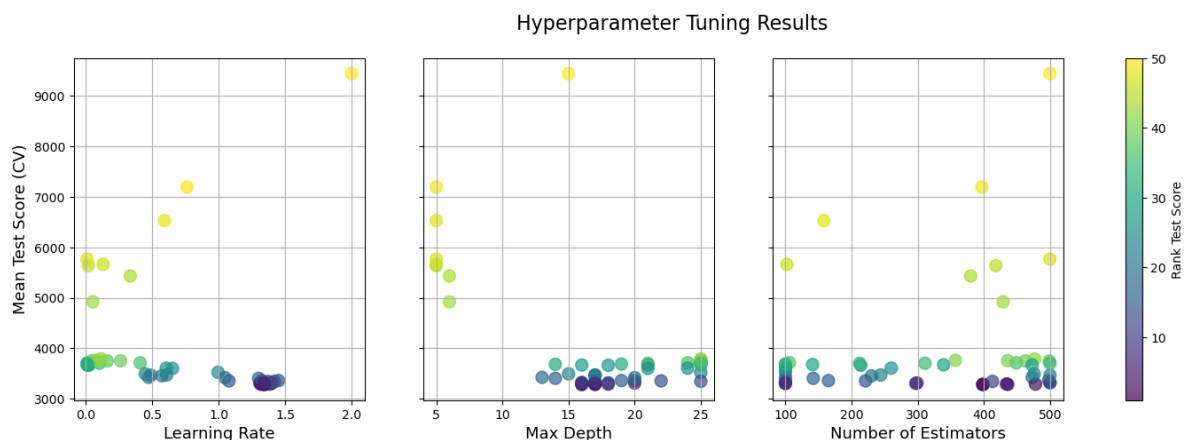
```
# Create 1x3 subplots
fig, axes = plt.subplots(1, 3, figsize=(18, 5), sharey=True)

# List of hyperparameters and axis labels
params = ['param_learning_rate', 'param_estimator__max_depth', 'param_n_estimators']
labels = ['Learning Rate', 'Max Depth', 'Number of Estimators']

# Plot each subplot
for ax, param, label in zip(axes, params, labels):
    sc = ax.scatter(
        results_df[param],
        results_df['mean_test_score'],
        c=results_df['rank_test_score'],
        cmap='viridis',
        s=100,
        alpha=0.7
    )
    ax.set_xlabel(label, fontsize=13)
    ax.grid(True)

# Set shared y-axis label and title
axes[0].set_ylabel('Mean Test Score (CV)', fontsize=13)
fig.suptitle('Hyperparameter Tuning Results', fontsize=16)

# Add shared colorbar
cbar = fig.colorbar(sc, ax=axes.ravel().tolist(), label='Rank Test Score');
```





### 3D scatterplot

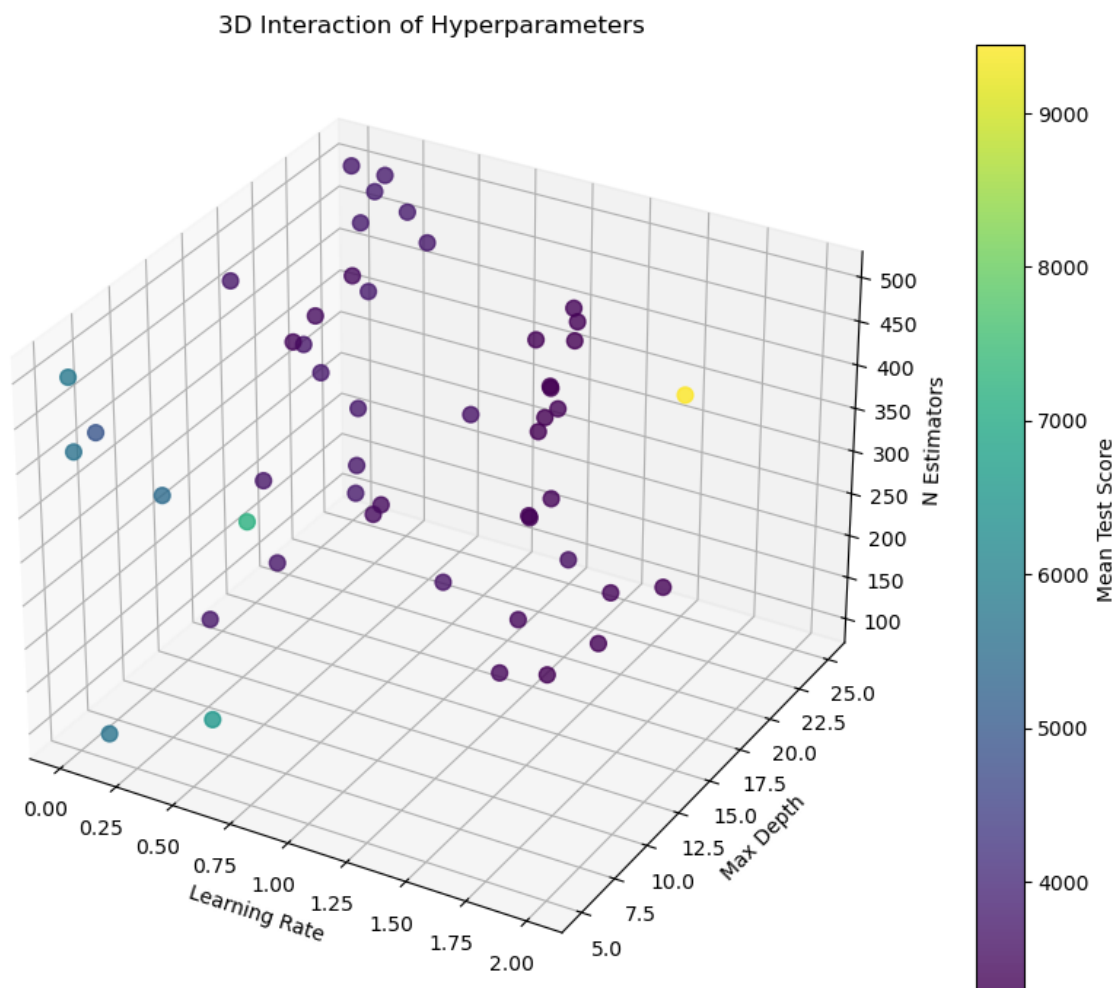
- Each point is a combination of the 3 hyperparameters.
- Color indicates performance (darker = better).
- You can rotate the 3D plot in Jupyter interactively!

```
from mpl_toolkits.mplot3d import Axes3D
import matplotlib.pyplot as plt

fig = plt.figure(figsize=(10, 7))
ax = fig.add_subplot(111, projection='3d')

p = ax.scatter(
    results_df['param_learning_rate'],
    results_df['param_estimator__max_depth'],
    results_df['param_n_estimators'],
    c=results_df['mean_test_score'],
    cmap='viridis',
    s=60,
    alpha=0.8
)

ax.set_xlabel('Learning Rate')
ax.set_ylabel('Max Depth')
ax.set_zlabel('N Estimators')
fig.colorbar(p, label='Mean Test Score')
plt.title('3D Interaction of Hyperparameters')
plt.tight_layout()
plt.show()
```



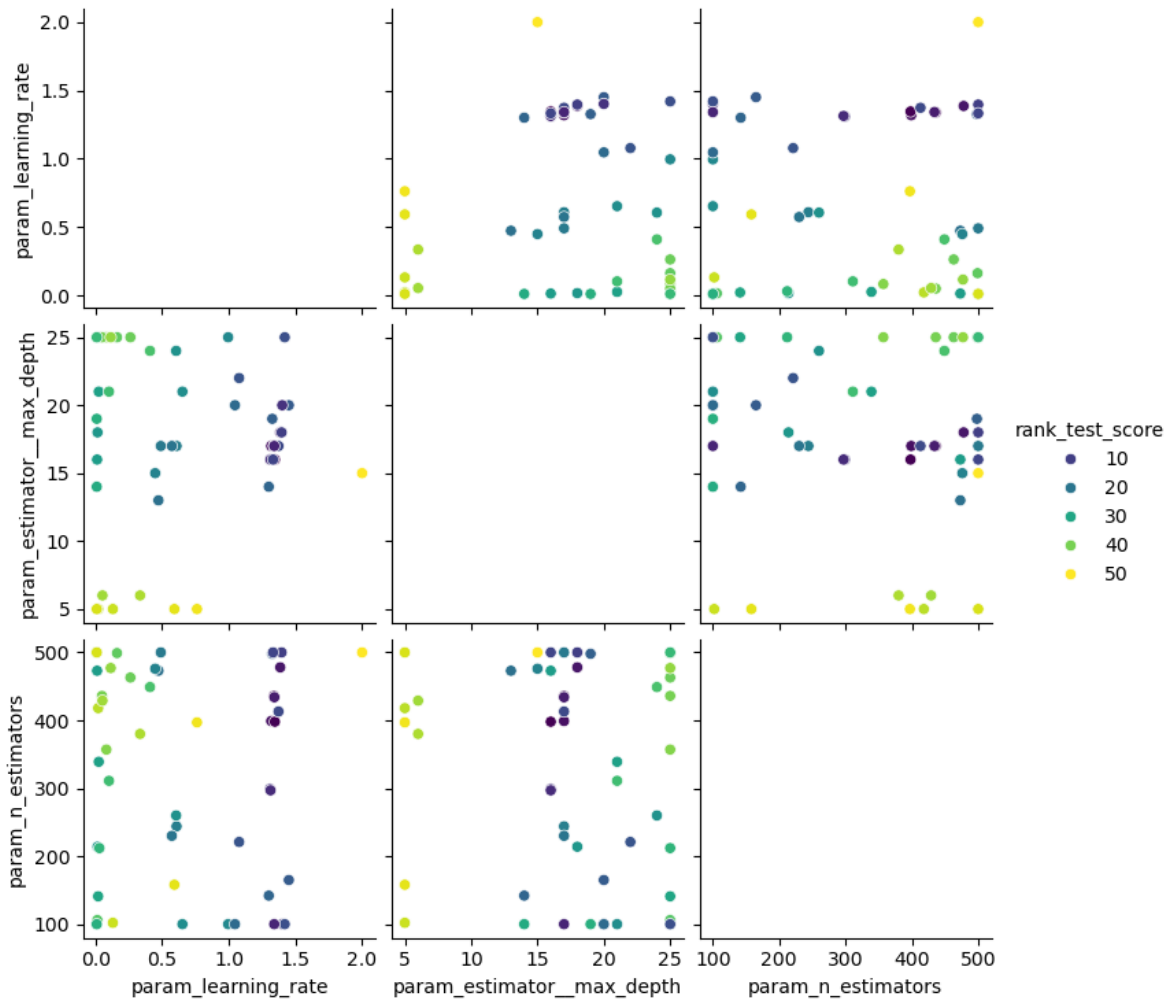
```
sorted_results = results_df.copy()
sorted_results = sorted_results[['param_learning_rate', 'param_estimator__max_depth', 'param_
sorted_results = sorted_results.sort_values(by='rank_test_score')
sorted_results.reset_index(drop=True, inplace=True)
sorted_results[:10] # Display the top 10 results
```

	param_learning_rate	param_estimator__max_depth	param_n_estimators	mean_test_score	std_test_score
0	1.346028	16	398	3280.273310	368.7
1	1.339271	17	436	3285.379694	347.7
2	1.318424	17	399	3286.757311	349.7
3	1.386113	18	478	3289.156528	363.7
4	1.342019	17	434	3291.786531	354.7

	param_learning_rate	param_estimator__max_depth	param_n_estimators	mean_test_score	std_test_score
5	1.340962	17	100	3300.807523	343.111111
6	1.313334	16	297	3304.849191	342.111111
7	1.401063	20	100	3309.832144	341.111111
8	1.309088	16	299	3312.005185	332.111111
9	1.396356	18	500	3317.350680	334.111111

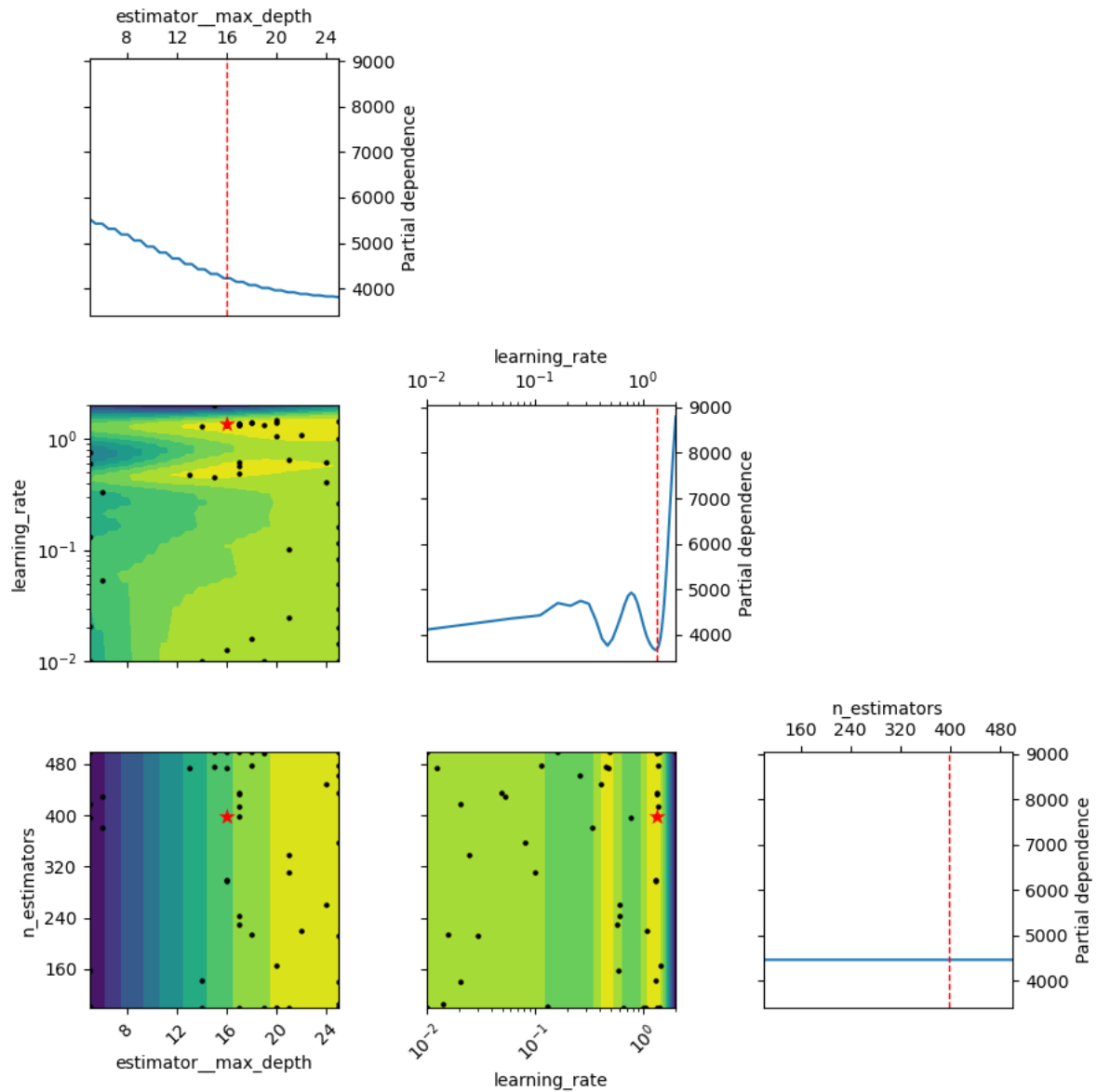
Let's analyze radeoffs/interactions

```
sns.pairplot(
    results_df,
    vars=[
        'param_learning_rate',
        'param_estimator__max_depth',
        'param_n_estimators'
    ],
    hue='rank_test_score',
    palette='viridis'
);
```



```
from skopt.plots import plot_objective

plot_objective(opt.optimizer_results_[0], dimensions=None, size = 3)
plt.tight_layout()
plt.show()
```



## 8.5.6 Using Optuna for Hyperparameter Tuning

```
pip install optuna
```

Collecting optuna

Downloading optuna-4.3.0-py3-none-any.whl.metadata (17 kB)

```

Collecting alembic>=1.5.0 (from optuna)
  Downloading alembic-1.15.2-py3-none-any.whl.metadata (7.3 kB)
Collecting colorlog (from optuna)
  Downloading colorlog-6.9.0-py3-none-any.whl.metadata (10 kB)
Requirement already satisfied: numpy in c:\users\lsi8012\appdata\local\anaconda3\lib\site-pa
Requirement already satisfied: packaging>=20.0 in c:\users\lsi8012\appdata\roaming\python\py
Requirement already satisfied: sqlalchemy>=1.4.2 in c:\users\lsi8012\appdata\local\anaconda3
Requirement already satisfied: tqdm in c:\users\lsi8012\appdata\local\anaconda3\lib\site-pac
Requirement already satisfied: PyYAML in c:\users\lsi8012\appdata\local\anaconda3\lib\site-p
Collecting Mako (from alembic>=1.5.0->optuna)
  Downloading mako-1.3.10-py3-none-any.whl.metadata (2.9 kB)
Collecting typing-extensions>=4.12 (from alembic>=1.5.0->optuna)
  Downloading typing_extensions-4.13.2-py3-none-any.whl.metadata (3.0 kB)
Requirement already satisfied: greenlet!=0.4.17 in c:\users\lsi8012\appdata\local\anaconda3\
Requirement already satisfied: colorama in c:\users\lsi8012\appdata\roaming\python\python312
Requirement already satisfied: MarkupSafe>=0.9.2 in c:\users\lsi8012\appdata\local\anaconda3
Downloading optuna-4.3.0-py3-none-any.whl (386 kB)
----- 0.0/386.6 kB ? eta -:--:--
----- 174.1/386.6 kB 3.5 MB/s eta 0:00:01
----- 386.6/386.6 kB 4.8 MB/s eta 0:00:00
Downloading alembic-1.15.2-py3-none-any.whl (231 kB)
----- 0.0/231.9 kB ? eta -:--:--
----- 231.9/231.9 kB 7.2 MB/s eta 0:00:00
Downloading colorlog-6.9.0-py3-none-any.whl (11 kB)
Downloading typing_extensions-4.13.2-py3-none-any.whl (45 kB)
----- 0.0/45.8 kB ? eta -:--:--
----- 45.8/45.8 kB 2.4 MB/s eta 0:00:00
Downloading mako-1.3.10-py3-none-any.whl (78 kB)
----- 0.0/78.5 kB ? eta -:--:--
----- 78.5/78.5 kB 4.6 MB/s eta 0:00:00
Installing collected packages: typing-extensions, Mako, colorlog, alembic, optuna
  Attempting uninstall: typing-extensions
    Found existing installation: typing_extensions 4.11.0
    Uninstalling typing_extensions-4.11.0:
      Successfully uninstalled typing_extensions-4.11.0
Successfully installed Mako-1.3.10 alembic-1.15.2 colorlog-6.9.0 optuna-4.3.0 typing-extensi
Note: you may need to restart the kernel to use updated packages.

```

ERROR: pip's dependency resolver does not currently take into account all the packages that a streamlit 1.32.0 requires packaging<24,>=16.8, but you have packaging 24.2 which is incompat.

Step 1: Import

```
# import optuna
import optuna
```

Step 2: Define the Objective Function

```
def objective(trial):
    # Suggest hyperparameters
    learning_rate = trial.suggest_float("learning_rate", 0.01, 2.0)
    max_depth = trial.suggest_int("max_depth", 5, 25)
    n_estimators = trial.suggest_int("n_estimators", 100, 500)

    # Define model with trial parameters
    base_estimator = DecisionTreeRegressor(max_depth=max_depth)
    model = AdaBoostRegressor(
        estimator=base_estimator,
        learning_rate=learning_rate,
        n_estimators=n_estimators,
        random_state=42
    )

    # Cross-validation score (negative RMSE)
    score = cross_val_score(model, X_train_final, y_train, scoring="neg_root_mean_squared_error")
    return -np.mean(score)
```

Step 3: Run the study

```
study = optuna.create_study(direction="minimize")
study.optimize(objective, n_trials=20, timeout=600) # 50 trials or 10 min
```

```
[I 2025-04-29 18:13:37,366] A new study created in memory with name: no-name-8c1fd49f-a877-4
[I 2025-04-29 18:13:55,695] Trial 0 finished with value: 3604.108244895723 and parameters: {
[I 2025-04-29 18:14:39,369] Trial 1 finished with value: 3418.989299711666 and parameters: {
[I 2025-04-29 18:15:01,876] Trial 2 finished with value: 6675.343723783454 and parameters: {
[I 2025-04-29 18:15:51,382] Trial 3 finished with value: 3534.3761450527854 and parameters: -
[I 2025-04-29 18:16:17,255] Trial 4 finished with value: 4146.601971284103 and parameters: {
[I 2025-04-29 18:16:51,329] Trial 5 finished with value: 3373.6473210617305 and parameters: -
[I 2025-04-29 18:17:22,661] Trial 6 finished with value: 4645.349739417665 and parameters: {
[I 2025-04-29 18:17:48,173] Trial 7 finished with value: 3783.098791065398 and parameters: {
[I 2025-04-29 18:18:07,203] Trial 8 finished with value: 3642.6233815254373 and parameters: -
[I 2025-04-29 18:19:28,631] Trial 9 finished with value: 3815.355472103091 and parameters: {
[I 2025-04-29 18:19:52,259] Trial 10 finished with value: 3460.264767659689 and parameters: -
```

```
[I 2025-04-29 18:20:24,557] Trial 11 finished with value: 3354.0174995460206 and parameters:
[I 2025-04-29 18:20:51,602] Trial 12 finished with value: 3424.5702876072514 and parameters:
[I 2025-04-29 18:21:42,335] Trial 13 finished with value: 3428.9633976083956 and parameters:
[I 2025-04-29 18:22:07,861] Trial 14 finished with value: 3363.7031695680153 and parameters:
[I 2025-04-29 18:22:38,820] Trial 15 finished with value: 3467.0742139036493 and parameters:
[I 2025-04-29 18:23:03,761] Trial 16 finished with value: 4339.523251732014 and parameters:
[I 2025-04-29 18:23:32,216] Trial 17 finished with value: 3429.3383627414196 and parameters:
[I 2025-04-29 18:23:46,880] Trial 18 finished with value: 3432.187531277779 and parameters:
```

#### Step 4: Review Best Result

```
print("Best RMSE:", study.best_value)
print("Best hyperparameters:", study.best_params)
```

Best RMSE: 3354.0174995460206

Best hyperparameters: {'learning\_rate': 1.4283231452601248, 'max\_depth': 17, 'n\_estimators':

```
# make a prediction using the best hyperparameters
best_params = study.best_params
base_estimator = DecisionTreeRegressor(max_depth=best_params['max_depth'])
model = AdaBoostRegressor(
    estimator=base_estimator,
    learning_rate=best_params['learning_rate'],
    n_estimators=best_params['n_estimators'],
    random_state=42
)
# fit the model
model.fit(X_train_final, y_train)
# predict the test set
y_pred = model.predict(X_test_final)
# calculate the mean squared error
rmse = root_mean_squared_error(y_test, y_pred)
print(f'RMSE: {rmse:.2f}')
print(f'R2 Score: {r2_score(y_test, y_pred)}')
```

RMSE: 3524.38

R2 Score: 0.9576611696103872

#### Step 5: Visualize the Search



```
optuna.visualization.plot_optimization_history(study).show()
```

Unable to display output for mime type(s): application/vnd.plotly.v1+json

```
optuna.visualization.plot_param_importances(study).show()
```

Unable to display output for mime type(s): application/vnd.plotly.v1+json

```
optuna.visualization.plot_parallel_coordinate(study).show()
```

Unable to display output for mime type(s): application/vnd.plotly.v1+json

### Insights:

- Best Hyperparameter Region:
  - `learning_rate`: ~0.2–0.5
  - `max_depth`: ~7–10
  - `n_estimators`: ~300–450
- Trade-offs:
  - Increasing `n_estimators` improves performance but increases computation time.
  - Lower `learning_rate` values require more estimators to achieve good performance.
- Next Steps:
  - Focus on fine-tuning within the identified ranges for `learning_rate`, `max_depth`, and `n_estimators`.
  - Use these insights to narrow the search space for further optimization.

## 9 Gradient Boosting

Check the gradient boosting algorithm in section 10.10.2 of the book, [Elements of Statistical Learning](#) before using these notes.

*Note that in this course, lecture notes are not sufficient, you must read the book for better understanding. Lecture notes are just implementing the concepts of the book on a dataset, but not explaining the concepts elaborately.*

### 9.1 Core Hyperparameters in Gradient Boosting

#### 1. Number of Trees (`n_estimators`)

- Use **early stopping** (via `n_iter_no_change` and `validation_fraction` in `scikit-learn`) to avoid overfitting.
- Start with a large value (e.g., 500–1000) and let early stopping prune unnecessary trees.

#### 2. Early Stopping

- Prevents overfitting by halting training once the validation performance stops improving.
- Controlled using:
  - `n_iter_no_change`: Number of rounds with no improvement before stopping (e.g., 10).
  - `validation_fraction`: Fraction of training data reserved as internal validation set (e.g., 0.1).
  - `tol`: Minimum improvement to be considered significant (e.g., `1e-4`).
- Set a large `n_estimators`, and let early stopping determine the optimal number of boosting iterations.

#### 3. Learning Rate (`learning_rate`)

- Shrinks the contribution of each tree to improve generalization.

- *Typical range:* 0.01–0.2 (lower values require more trees).

#### 4. Tree Complexity

- `max_depth`: Depth of individual trees. Start shallow (3–6) to limit overfitting.
- `min_samples_split`: Minimum samples required to split a node (e.g., 10–50).
- `min_samples_leaf`: Minimum samples required in a leaf node (e.g., 5–20).

#### 5. Stochastic Gradient Boosting

- `subsample`: Fraction of training data sampled per tree (e.g., 0.5–1.0).
- `max_features`: Fraction/absolute number of features used per split (e.g., `sqrt(n_features)` or 0.8).

#### 6. Loss Function (loss)

- Matches the problem type:
  - Regression: `squared_error`, `absolute_error`
  - Classification: `log_loss` (binary/multinomial deviance)

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import cross_val_score, train_test_split, KFold, cross_val_predict
from sklearn.metrics import root_mean_squared_error, mean_squared_error, r2_score, roc_curve, auc,
recall_score, precision_score, confusion_matrix
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV, ParameterGrid, StratifiedKFold
from sklearn.ensemble import GradientBoostingRegressor, GradientBoostingClassifier, BaggingRegressor
from sklearn.preprocessing import OneHotEncoder, FunctionTransformer
import itertools as it
import time as time

from skopt import BayesSearchCV
from skopt.space import Real, Categorical, Integer
from skopt.plots import plot_objective, plot_histogram, plot_convergence
import warnings
from IPython import display
```

Let's reuse the car dataset to evaluate how different hyperparameter settings affect the performance of gradient boosting

```
# Load the dataset
car = pd.read_csv('Datasets/car.csv')
car.head()
```

	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	vw	Beetle	2014	Manual	55457	Diesel	30	65.3266	1.6	7490
1	vauxhall	GTC	2017	Manual	15630	Petrol	145	47.2049	1.4	10998
2	merc	G Class	2012	Automatic	43000	Diesel	570	25.1172	3.0	44990
3	audi	RS5	2019	Automatic	10	Petrol	145	30.5593	2.9	51990
4	merc	X-CLASS	2018	Automatic	14000	Diesel	240	35.7168	2.3	28990

```
X = car.drop(columns=['price'])
y = car['price']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# extract the categorical columns and put them in a list
categorical_feature = X.select_dtypes(include=['object']).columns.tolist()

# extract the numerical columns and put them in a list
numerical_feature = X.select_dtypes(include=['int64', 'float64']).columns.tolist()

encoder = OneHotEncoder(handle_unknown='ignore', sparse_output=False)

X_train_encoded = encoder.fit_transform(X_train[categorical_feature])
X_test_encoded = encoder.transform(X_test[categorical_feature])

# Convert the encoded features back to DataFrame
X_train_encoded_df = pd.DataFrame(X_train_encoded, columns=encoder.get_feature_names_out(categorical_feature))
X_test_encoded_df = pd.DataFrame(X_test_encoded, columns=encoder.get_feature_names_out(categorical_feature))

# Concatenate the encoded features with the original numerical features
X_train_final = pd.concat([X_train_encoded_df, X_train[numerical_feature].reset_index(drop=True)], axis=1)
X_test_final = pd.concat([X_test_encoded_df, X_test[numerical_feature].reset_index(drop=True)], axis=1)
```

### 9.1.1 Effect of Number of Trees on Cross-Validation Error

Effect of Number of Trees on Cross-Validation Error In Gradient Boosting, the number of trees (`n_estimators`) controls how many boosting rounds the model performs. Adding more trees can reduce bias and improve training accuracy, but it also increases the risk of overfitting, especially with a high learning rate.

The optimal number of trees is often found by balancing **model complexity** and **generalization performance** using cross-validation.

```
def get_models():
    models = dict()
    # define number of trees to consider
    n_trees = [2, 5, 10, 50, 100, 500, 1000, 2000, 5000]
    for n in n_trees:
        models[str(n)] = GradientBoostingRegressor(n_estimators=n, random_state=1, loss='huber')
    return models

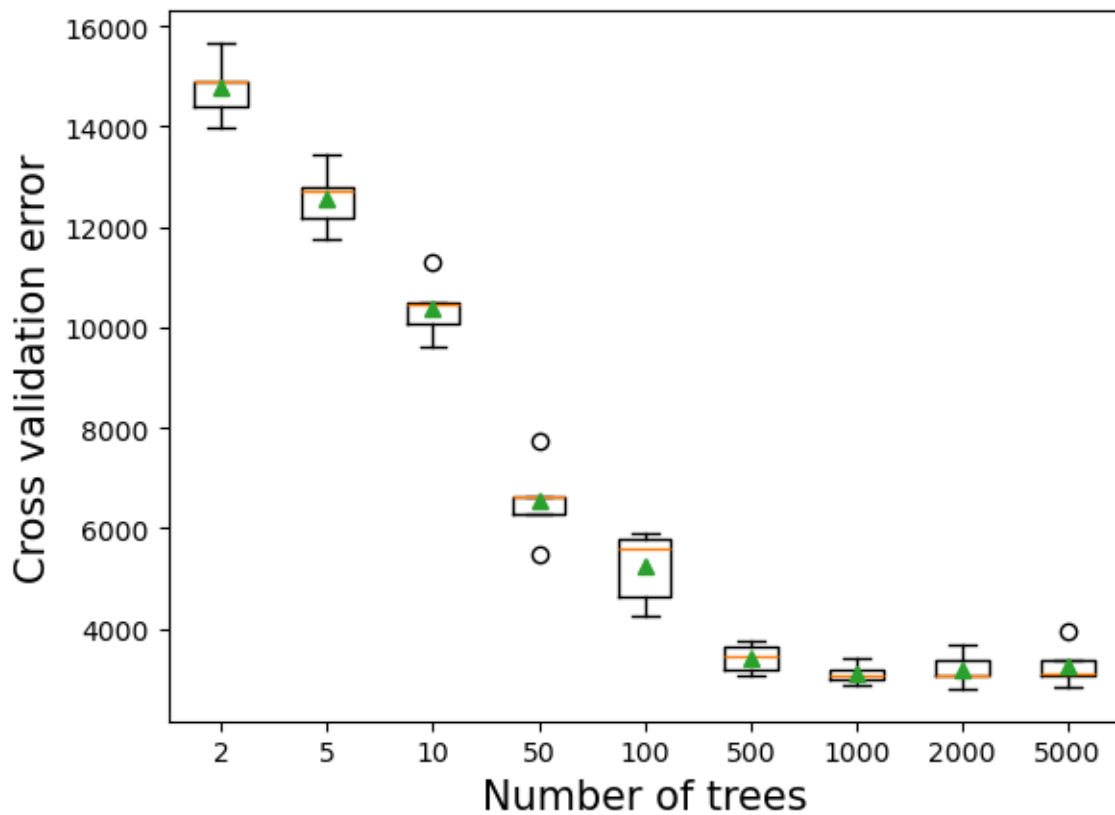
# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=5, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = np.sqrt(-cross_val_score(model, X, y, scoring='neg_mean_squared_error', cv=cv, n_jobs=-1))
    return scores

# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X_train_final, y_train)
    # store the results
    results.append(scores)
    names.append(name)
    # summarize the performance along the way
    print('>%s %.3f (%.3f)' % (name, np.mean(scores), np.std(scores)))
# plot model performance for comparison
plt.boxplot(results, labels=names, showmeans=True)
plt.ylabel('Cross validation error', fontsize=15)
plt.xlabel('Number of trees', fontsize=15);
```

```

>2 14769.612 (563.444)
>5 12572.047 (567.992)
>10 10383.351 (555.641)
>50 6549.576 (722.462)
>100 5232.949 (656.216)
>500 3419.467 (262.753)
>1000 3106.002 (184.799)
>2000 3194.874 (293.134)
>5000 3272.374 (380.296)

```



### 9.1.2 Early stopping in Gradient Boosting

**Early stopping** is a regularization technique used to prevent overfitting in boosting models. Instead of specifying a fixed number of trees (`n_estimators`), the algorithm monitors performance on a **validation set** and stops adding new trees once the model's improvement has plateaued.

In scikit-learn, early stopping can be enabled using: - `early_stopping=True` - `validation_fraction`: The fraction of training data used as a validation set - `n_iter_no_change`: Number of iterations to wait without improvement before stopping

This approach not only improves generalization but also reduces training time by avoiding unnecessary trees.

```
params = dict(n_estimators=2000, max_depth=5, learning_rate=0.1, random_state=42)

gbm_full = GradientBoostingRegressor(**params)
gbm_early_stopping = GradientBoostingRegressor(
    **params,
    validation_fraction=0.1,
    n_iter_no_change=10,
)

start_time = time.time()
gbm_full.fit(X_train_final, y_train)
training_time_full = time.time() - start_time
n_estimators_full = gbm_full.n_estimators_

start_time = time.time()
gbm_early_stopping.fit(X_train_final, y_train)
training_time_early_stopping = time.time() - start_time
estimators_early_stopping = gbm_early_stopping.n_estimators_
```

Let's calculate the RMSE on both the training and test datasets for each model, which will be used for later visualization.

```
# import root mean squared error function
from sklearn.metrics import root_mean_squared_error

train_errors_without = []
test_errors_without = []

train_errors_with = []
test_errors_with = []

for i, (train_pred, test_pred) in enumerate(
    zip(
        gbm_full.staged_predict(X_train_final),
        gbm_full.staged_predict(X_test_final),
    )
)
```

```

):
    train_errors_without.append(root_mean_squared_error(y_train, train_pred))
    test_errors_without.append(root_mean_squared_error(y_test, test_pred))

for i, (train_pred, test_pred) in enumerate(
    zip(
        gbm_early_stopping.staged_predict(X_train_final),
        gbm_early_stopping.staged_predict(X_test_final),
    )
):
    train_errors_with.append(root_mean_squared_error(y_train, train_pred))
    test_errors_with.append(root_mean_squared_error(y_test, test_pred))

```

Let's visualize Comparison. It includes three subplots:

1. Plotting training errors of both models over boosting iterations.
2. Plotting test errors of both models over boosting iterations.
3. Creating a bar chart to compare the training times and the number of estimators used by the models with and without early stopping.

```

fig, axes = plt.subplots(ncols=3, figsize=(12, 4))

axes[0].plot(train_errors_without, label="gbm_full")
axes[0].plot(train_errors_with, label="gbm_early_stopping")
axes[0].set_xlabel("Boosting Iterations")
axes[0].set_ylabel("RMSE (Training)")
axes[0].set_yscale("log")
axes[0].legend()
axes[0].set_title("Training Error")

axes[1].plot(test_errors_without, label="gbm_full")
axes[1].plot(test_errors_with, label="gbm_early_stopping")
axes[1].set_xlabel("Boosting Iterations")
axes[1].set_ylabel("RMSE (Test)")
axes[1].set_yscale("log")
axes[1].legend()
axes[1].set_title("Test Error")

training_times = [training_time_full, training_time_early_stopping]
labels = ["gbm_full", "gbm_early_stopping"]
bars = axes[2].bar(labels, training_times)
axes[2].set_ylabel("Training Time (s)")

```

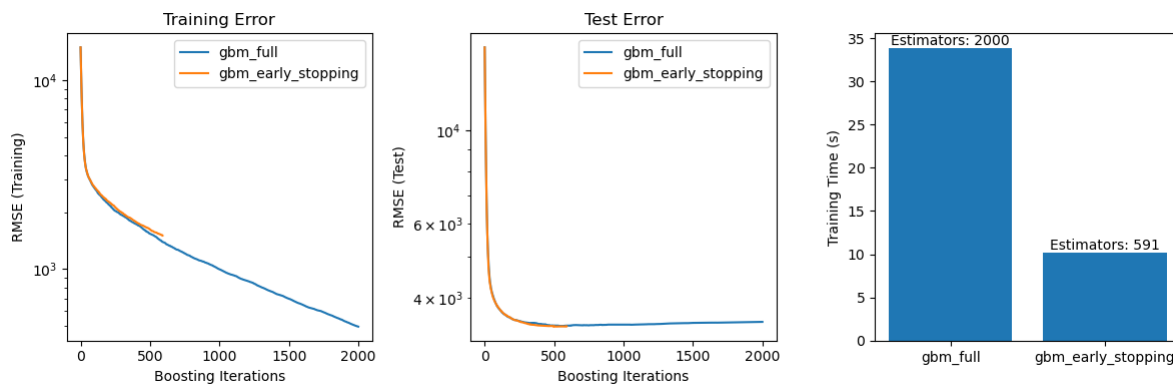


```

for bar, n_estimators in zip(bars, [n_estimators_full, estimators_early_stopping]):
    height = bar.get_height()
    axes[2].text(
        bar.get_x() + bar.get_width() / 2,
        height + 0.001,
        f"Estimators: {n_estimators}",
        ha="center",
        va="bottom",
    )

plt.tight_layout()
plt.show()

```



The difference in training error between the `gbm_full` and the `gbm_early_stopping` stems from the fact that `gbm_early_stopping` sets aside `validation_fraction` of the training data as an internal validation set.

Early stopping is decided based on this internal validation score.

### Benefits of Using Early Stopping in Boosting:

- **Preventing Overfitting**

Early stopping helps avoid overfitting by monitoring the test error.

When the error stabilizes or starts increasing, training stops — resulting in better generalization to unseen data.

- **Improving Training Efficiency**

Models with early stopping often require **fewer estimators** while achieving similar accuracy.

This reduces training time significantly compared to training without early stopping.

### 9.1.3 Effect of Learning Rate on Cross-Validation Error

The learning rate (`learning_rate`) determines how much each new tree contributes to the overall model. A **smaller learning rate** results in slower learning and often requires more trees to achieve good performance. A **larger learning rate** speeds up learning but increases the risk of overfitting.

Finding the optimal learning rate involves balancing: - **High learning rate** → faster convergence, but higher risk of overfitting

- **Low learning rate** → better generalization, but requires more trees and longer training time

Cross-validation helps identify the learning rate that minimizes prediction error while ensuring model stability.

```
def get_models():
    models = dict()
    # create 9 evenly spaced values between 0.2 and 1.0
    learning_rates = np.linspace(0.2, 1.0, 9)
    for learning_rate in learning_rates:
        # Round to 2 decimal places for clean keys
        lr_rounded = round(learning_rate, 2)
        key = f"{lr_rounded:.2f}"
        models[key] = GradientBoostingRegressor(learning_rate=lr_rounded, random_state=1, loss='ls')
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=5, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = np.sqrt(-cross_val_score(model, X, y, scoring='neg_mean_squared_error', cv=cv, n_jobs=-1))
    return scores

# get the models to evaluate
models = get_models()

# evaluate the models and store results
results, names = list(), list()
mean_scores = [] # Track mean scores separately

for name, model in models.items():
    # evaluate the model
```

```

    scores = evaluate_model(model, X_train_final, y_train)
    # store the results
    results.append(scores)
    names.append(name)
    # Calculate and store mean score
    mean_score = np.mean(scores)
    mean_scores.append(mean_score)
    # summarize the performance along the way
    print('>%s %.1f (%.1f)' % (name, mean_score, np.std(scores)))

# plot model performance for comparison
plt.figure(figsize=(10, 7))
plt.boxplot(results, labels=names, showmeans=True)
plt.ylabel('Cross validation error', fontsize=15)
plt.xlabel('Learning rate', fontsize=15)
plt.title('Model Performance by Learning Rate', fontsize=16)
plt.grid(True, linestyle='--', alpha=0.7)

# Find the best model using the saved mean scores
best_index = np.argmin(mean_scores)
best_lr = names[best_index]
best_score = mean_scores[best_index]

# Highlight the best model on the plot
plt.axvline(x=best_index+1, color='red', linestyle='--', alpha=0.7)
plt.text(best_index+1.2, min(mean_scores)*0.95,
         f'Best: {best_lr} (RMSE: {best_score:.1f})',
         color='red', fontweight='bold')

plt.show()

# Print the best model information
print(f"\nBest model: {best_lr} with RMSE: {best_score:.3f}")

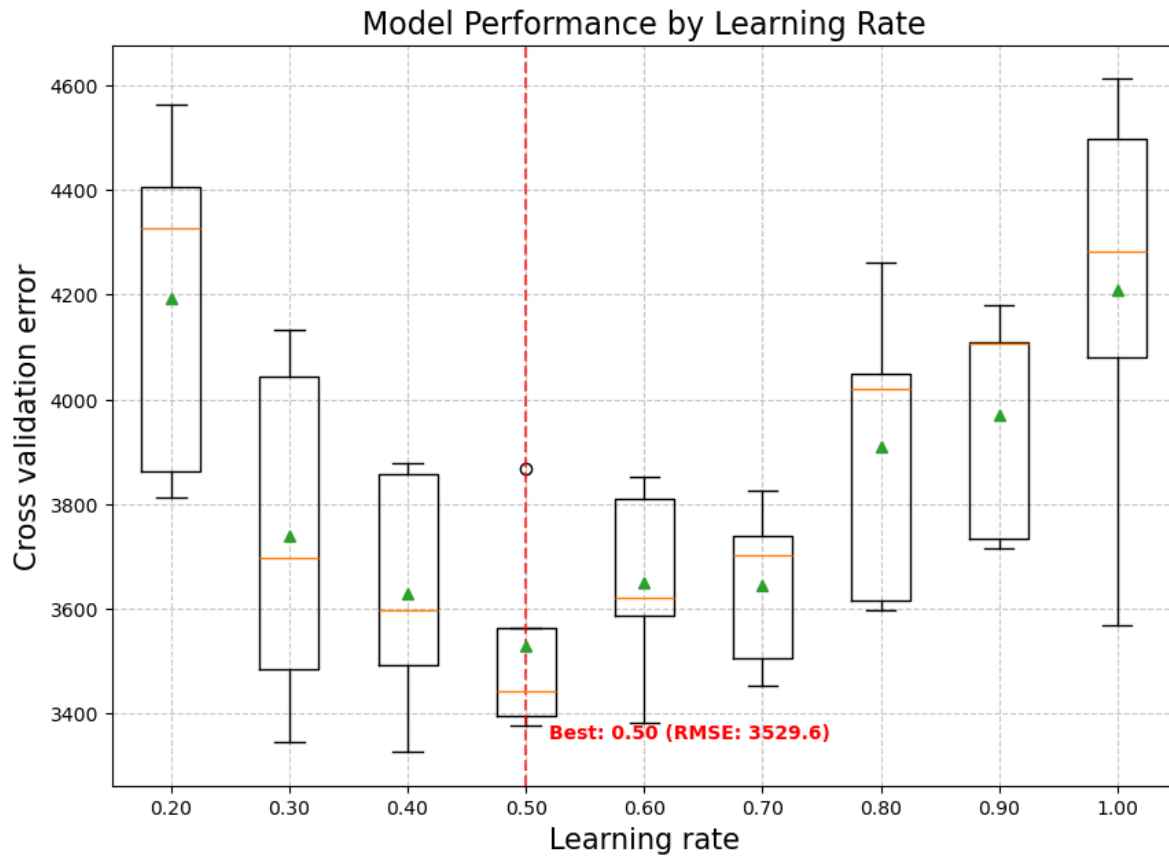
```

```

>0.20 4193.7 (301.2)
>0.30 3740.3 (306.3)
>0.40 3630.0 (212.2)
>0.50 3529.6 (181.5)
>0.60 3650.2 (169.0)
>0.70 3644.7 (142.5)
>0.80 3908.9 (260.6)
>0.90 3968.7 (201.1)

```

>1.00 4208.3 (368.5)



Best model: 0.50 with RMSE: 3529.563

#### 9.1.4 Learning Rate and Number of Trees Are Closely Linked

The **learning rate** and **number of trees** (`n_estimators`) are tightly coupled hyperparameters in gradient boosting. Their balance plays a key role in model performance and overfitting control.

- A **lower learning rate** slows the learning process, requiring **more trees** to achieve strong performance.
- A **higher learning rate** speeds up training but may cause the model to **overfit** if not regularized properly.

A high learning rate with too few trees can lead to poor generalization, while a very low learning rate with too many trees may improve accuracy but increase training time significantly.

**Best practice:** Use a low to moderate learning rate (e.g., 0.01–0.1) combined with **early stopping** to find the optimal number of trees.

### 9.1.5 Effect of Stochastic Gradient Boosting on Cross-Validation Error

**Stochastic Gradient Boosting** enhances generalization by introducing randomness into the model-building process. Two key hyperparameters that control this are **subsample** and **max\_features**, and they operate on **different dimensions** of the data:

Parameter	Applies To	Purpose
<b>subsample</b>	Rows (data points)	Randomly samples a fraction of the training data for each tree
<b>max_features</b>	Columns (features)	Randomly samples a fraction of the features for each tree or split

By tuning these parameters, we can reduce overfitting and increase model robustness. However, setting them too low may lead to underfitting due to insufficient information per tree.

```
from sklearn.metrics import make_scorer, mean_squared_error

# Define model
model = GradientBoostingRegressor(n_estimators=100, max_depth=4, learning_rate=0.1, random_s

# Define param grid
param_grid = {
    'subsample': np.linspace(0.2, 1.0, 9),
    'max_features': np.linspace(0.2, 1.0, 9)
}

# RMSE scoring
scorer = make_scorer(mean_squared_error, greater_is_better=False)

# Grid search
grid = GridSearchCV(estimator=model, param_grid=param_grid,
                    scoring=scorer, cv=5, n_jobs=-1, verbose=1)
grid.fit(X_train_final, y_train)

# Create DataFrame from results
```

```
results_df = pd.DataFrame(grid.cv_results_)
results_df['mean_rmse'] = np.sqrt(-results_df['mean_test_score'])
```

Fitting 5 folds for each of 81 candidates, totalling 405 fits

```
# Round subsample and max_features to 2 decimal places for display
results_df['subsample'] = results_df['param_subsample'].astype(float).round(2)
results_df['max_features'] = results_df['param_max_features'].astype(float).round(2)

# Then pivot using the rounded values
heatmap_data = results_df.pivot(index='subsample', columns='max_features', values='mean_rmse')

# Plot heatmap
plt.figure(figsize=(12, 9))
sns.heatmap(heatmap_data, annot=True, fmt=".3f", cmap="YlGnBu", cbar_kws={'label': 'CV RMSE'})
plt.title('Grid Search: CV RMSE by Subsample and Max Features')
plt.ylabel('Subsample')
plt.xlabel('Max Features')
plt.tight_layout()
plt.show()
```



### 9.1.6 Effect of Tree Complexity on Cross-Validation Error (Not Tuned Here)

**Tree complexity** controls how expressive and flexible each individual tree in the gradient boosting ensemble can be. While deeper and more complex trees can capture intricate patterns in the data, they are also more prone to overfitting, especially when combined with many trees.

Key parameters include:

- **max\_depth**: Limits the depth of each tree. Shallower trees (e.g., depth 3–6) are preferred for reducing overfitting.
- **min\_samples\_split**: Specifies the minimum number of samples required to split an internal node. Higher values make the tree more conservative.
- **min\_samples\_leaf**: Sets the minimum number of samples required to be at a leaf node. This also helps smooth the model and avoid capturing noise.

These parameters influence the bias-variance trade-off by adjusting how expressive each tree can be.

Since we have already discussed and tuned these parameters in earlier lessons (decision trees and random forests), we will **not tune them again here**.

### 9.1.7 Loss Function (`loss`)

In gradient boosting, the loss function determines how the model measures prediction errors and guides the optimization process during training. Here's a breakdown of common loss functions for regression and classification tasks:

- **Regression:**
  - `squared_error`: Penalizes larger errors more heavily; sensitive to outliers. (*Default for regression*)
  - `absolute_error`: Penalizes all errors equally; more robust to outliers.
  - `huber`: Combines squared and absolute error; less sensitive to outliers than `squared_error` and smoother than `absolute_error`.
- **Classification:**
  - `log_loss`: Also known as logistic loss or deviance; commonly used for binary and multiclass classification.
  - `exponential`: Used by AdaBoost; heavily penalizes misclassified points, making it more sensitive to outliers.

Choosing an appropriate loss function ensures the model is optimized for the specific structure and goals of the problem.

## 9.2 Joint Hyperparameter Tuning in Gradient Boosting for Regression

Since the optimal values of hyperparameters are often interdependent, they should be tuned **together** rather than in isolation to achieve the best performance.

### 9.2.1 Using `BayesSearchCV`

We can use `BayesSearchCV` with early stopping to **simultaneously tune multiple hyperparameters** in a more efficient and automated way.



```

# time the search
start = time.time()
# Define the search space
search_space = {
    'learning_rate': Real(0.01, 0.8, prior='log-uniform'), # Prefer lower rates
    'max_depth': Integer(4, 32), # Shallow trees to prevent overfitting
    'min_samples_split': Integer(2, 100), # Regularize splits
    'min_samples_leaf': Integer(1, 30), # Regularize leaves
    'subsample': Real(0.1, 1.0), # Stochastic sampling
    'max_features': Categorical([
        'sqrt', 'log2', None, # String options
        0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9 # Fractional options (discrete)
    ]) # Feature sampling
}

# Define the model
model_with_early_stopping = GradientBoostingRegressor(
    n_estimators=10000, # Start with a large number of trees
    validation_fraction=0.1, # Reserve 10% of training data for validation
    n_iter_no_change=10, # Stop after 20 rounds of no improvement
    tol=0.001, # Tolerance for early stopping
    random_state=42
)

# Define the search
bayes_cv = BayesSearchCV(
    model_with_early_stopping,
    search_space,
    n_iter=50, # Number of iterations
    scoring='neg_mean_squared_error',
    cv=5, # Cross-validation folds
    n_jobs=-1, # Use all available cores
    verbose=1, # Verbosity level
    random_state=42 # For reproducibility
)

# Fit the model
bayes_cv.fit(X_train_final, y_train)
# Stop the timer
end = time.time()
# Calculate elapsed time
elapsed_time = (end - start)/60 # Convert to minutes
# Print elapsed time
print(f"Elapsed time for Bayesian optimization with early stopping: {elapsed_time:.2f} minutes")

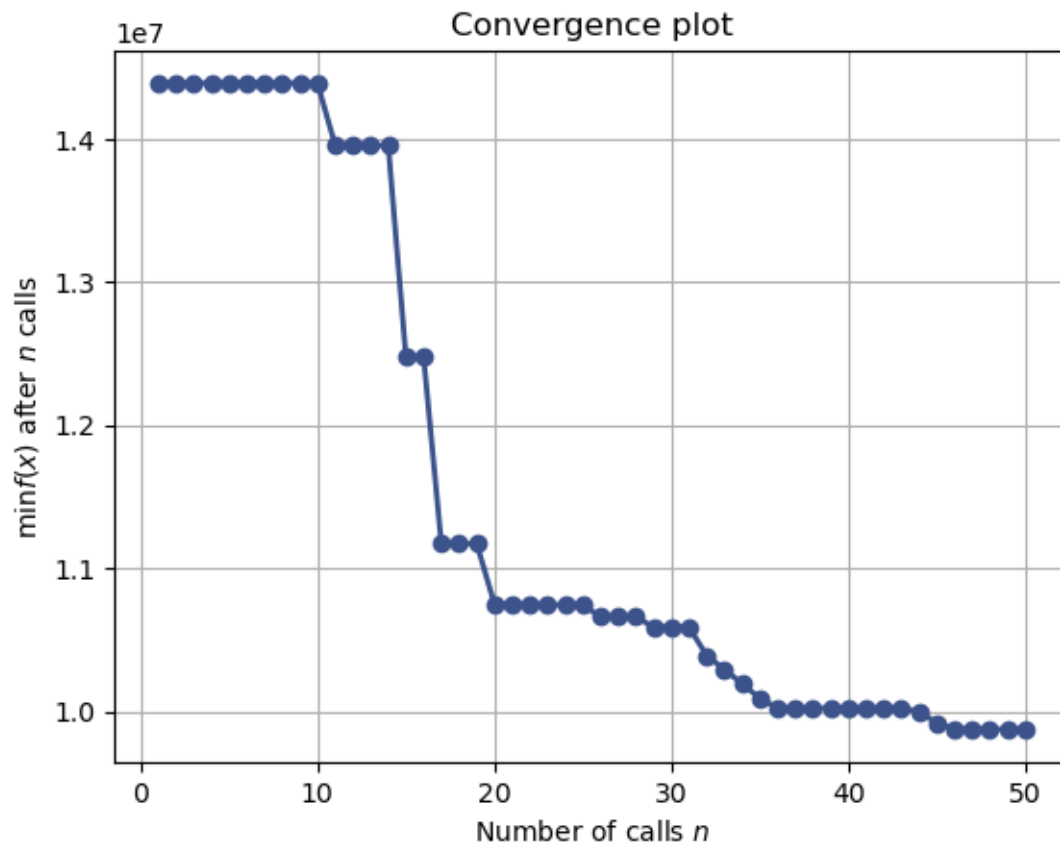
```

```
# Extract the best parameters and score
best_params = bayes_cv.best_params_
best_score = np.sqrt(-bayes_cv.best_score_)

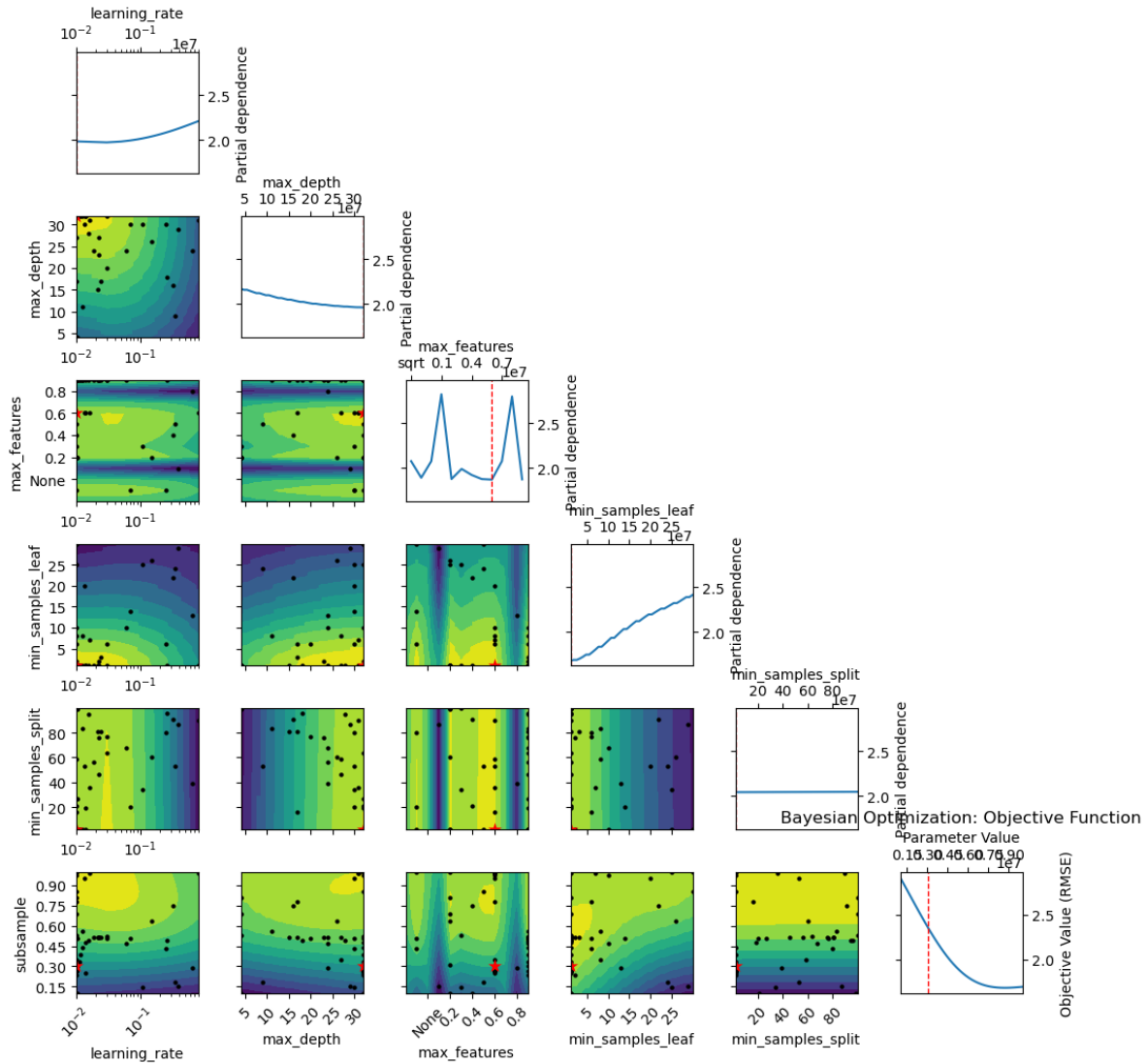
print(f"Best Parameters: {best_params}")
print(f"Best CV RMSE: {best_score:.3f}")
```

[illegible]





```
# Plot the objective function
plot_objective(bayes_cv.optimizer_results_[0])
plt.title('Bayesian Optimization: Objective Function')
plt.xlabel('Parameter Value')
plt.ylabel('Objective Value (RMSE)')
plt.show()
```



## 9.2.2 Tuning with Optuna

```
def objective(trial):
    # Define hyperparameters to optimize
    params = {
        'learning_rate': trial.suggest_float('learning_rate', 0.01, 0.8, log=True),
        'max_depth': trial.suggest_int('max_depth', 4, 32),
        'min_samples_split': trial.suggest_int('min_samples_split', 2, 100),
        'min_samples_leaf': trial.suggest_int('min_samples_leaf', 1, 30),
```

```

        'subsample': trial.suggest_float('subsample', 0.1, 1.0),
        'max_features': trial.suggest_categorical(
            'max_features',
            ['sqrt', 'log2', None, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]
        ),
        'n_iter_no_change': 10, # Stop if no improvement in 50 rounds
        'validation_fraction': 0.1, # 10% of training data for validation
        'tol': 0.001, # Tolerance for early stopping
        'n_estimators': 10000, # Start with a large number of trees
        'random_state': 42
    }

    # Initialize the model with the parameters, adding early stopping
    model = GradientBoostingRegressor(
        **params
    )
    model = GradientBoostingRegressor(**params)
    # Define the evaluation procedure
    cv = KFold(n_splits=5, shuffle=True, random_state=1)
    # Perform cross-validation
    scores = cross_val_score(model, X_train_final, y_train, cv=cv, scoring='neg_mean_squared_error')
    return np.mean(np.sqrt(-scores))

```

```

import optuna
start = time.time()
# Create a study object
study = optuna.create_study(direction="minimize")
study.optimize(objective, n_trials=50, timeout=600) # 50 trials or 10 min
# Stop the timer
end = time.time()
# Calculate elapsed time
elapsed_time = (end - start)/60 # Convert to minutes
print(f"Elapsed time for Optuna optimization: {elapsed_time:.2f} minutes")
# Extract the best parameters and score
best_params_optuna = study.best_params
best_score_optuna = study.best_value
print(f"Best Parameters: {best_params_optuna}")
print(f"Best CV RMSE: {best_score_optuna:.3f}")

```

```

[I 2025-05-06 17:02:10,428] A new study created in memory with name: no-name-a137b8c9-3b60-4...
[I 2025-05-06 17:02:29,091] Trial 0 finished with value: 3713.7523573640283 and parameters: -

```

[I 2025-05-06 17:02:30,785] Trial 1 finished with value: 4709.295547981693 and parameters: {

[I 2025-05-06 17:02:31,180] Trial 2 finished with value: 6213.951571057093 and parameters: {

[I 2025-05-06 17:02:31,883] Trial 3 finished with value: 6581.2879973994595 and parameters: {

[I 2025-05-06 17:02:32,144] Trial 4 finished with value: 5363.950030465088 and parameters: {

[I 2025-05-06 17:02:32,430] Trial 5 finished with value: 4845.887213572036 and parameters: {

[I 2025-05-06 17:02:33,509] Trial 6 finished with value: 4571.976448159823 and parameters: {

[I 2025-05-06 17:02:37,056] Trial 7 finished with value: 3877.957199433686 and parameters: {

[I 2025-05-06 17:02:43,003] Trial 8 finished with value: 4569.273660258111 and parameters: {

[I 2025-05-06 17:02:45,166] Trial 9 finished with value: 3698.820690997482 and parameters: {

[I 2025-05-06 17:02:51,446] Trial 10 finished with value: 3732.1319915304725 and parameters: {

[I 2025-05-06 17:02:58,569] Trial 11 finished with value: 4300.5959486610955 and parameters: {

[I 2025-05-06 17:02:59,894] Trial 12 finished with value: 4520.923335421491 and parameters: {

[I 2025-05-06 17:03:01,948] Trial 13 finished with value: 3571.01788037487 and parameters: {

[I 2025-05-06 17:03:04,511] Trial 14 finished with value: 3721.8669565390346 and parameters: {

[I 2025-05-06 17:03:08,463] Trial 15 finished with value: 4586.350862990998 and parameters: {

[I 2025-05-06 17:03:09,014] Trial 16 finished with value: 4058.9829855282683 and parameters: {

[I 2025-05-06 17:03:10,412] Trial 17 finished with value: 3517.5610120729484 and parameters: {

[I 2025-05-06 17:03:12,200] Trial 18 finished with value: 3430.5907640777577 and parameters: {

[I 2025-05-06 17:03:13,771] Trial 19 finished with value: 3799.704577632151 and parameters: {

[I 2025-05-06 17:03:16,441] Trial 20 finished with value: 3462.9715057349613 and parameters: {

[I 2025-05-06 17:03:19,025] Trial 21 finished with value: 3476.7550517115924 and parameters: {

[I 2025-05-06 17:03:21,047] Trial 22 finished with value: 3383.380528420081 and parameters: {

[I 2025-05-06 17:03:23,960] Trial 23 finished with value: 3209.6438264561025 and parameters: {

[I 2025-05-06 17:03:26,667] Trial 24 finished with value: 3317.3765168790715 and parameters: {

[I 2025-05-06 17:03:28,696] Trial 25 finished with value: 3288.8335695654087 and parameters: {

[I 2025-05-06 17:03:32,953] Trial 26 finished with value: 3425.413176428209 and parameters: {

[I 2025-05-06 17:03:35,890] Trial 27 finished with value: 3469.518619833793 and parameters: {

[I 2025-05-06 17:03:41,256] Trial 28 finished with value: 3555.0715331503548 and parameters: {

[I 2025-05-06 17:03:47,727] Trial 29 finished with value: 3257.3393769477807 and parameters: {

[I 2025-05-06 17:03:59,875] Trial 30 finished with value: 3455.808400573129 and parameters: {

[I 2025-05-06 17:04:01,814] Trial 31 finished with value: 3348.559154847241 and parameters: {

[I 2025-05-06 17:04:06,672] Trial 32 finished with value: 3323.7458320560645 and parameters: {

[I 2025-05-06 17:04:13,437] Trial 33 finished with value: 3424.053952590986 and parameters: {

[I 2025-05-06 17:04:16,575] Trial 34 finished with value: 3247.244108372873 and parameters: {

[I 2025-05-06 17:04:33,247] Trial 35 finished with value: 3512.98378266172 and parameters: {

[I 2025-05-06 17:04:42,133] Trial 36 finished with value: 3273.2818821200813 and parameters: {

[I 2025-05-06 17:04:59,100] Trial 37 finished with value: 3478.727227894916 and parameters: {

[I 2025-05-06 17:05:09,392] Trial 38 finished with value: 3309.6797560000705 and parameters: {

[I 2025-05-06 17:06:01,821] Trial 39 finished with value: 3755.080833012932 and parameters: {

[I 2025-05-06 17:06:11,332] Trial 40 finished with value: 3733.741397844383 and parameters: {

[I 2025-05-06 17:06:15,992] Trial 41 finished with value: 3306.4308993805207 and parameters: {

[I 2025-05-06 17:06:19,531] Trial 42 finished with value: 3215.3040996921914 and parameters: {

[I 2025-05-06 17:06:23,455] Trial 43 finished with value: 3234.8110020390777 and parameters: {

```
[I 2025-05-06 17:06:27,847] Trial 44 finished with value: 3249.4922397265354 and parameters:
[I 2025-05-06 17:06:31,136] Trial 45 finished with value: 3210.997902617934 and parameters:
[I 2025-05-06 17:06:34,258] Trial 46 finished with value: 3244.661570113128 and parameters:
[I 2025-05-06 17:06:37,225] Trial 47 finished with value: 3205.3588096764424 and parameters:
[I 2025-05-06 17:06:40,411] Trial 48 finished with value: 3245.293742974289 and parameters:
[I 2025-05-06 17:06:43,714] Trial 49 finished with value: 3238.761898699402 and parameters:
```

Elapsed time for Optuna optimization: 4.55 minutes

Best Parameters: {'learning\_rate': 0.022878117048627775, 'max\_depth': 23, 'min\_samples\_split': 2}

Best CV RMSE: 3205.359

```
optuna.visualization.plot_optimization_history(study).show()
```

Unable to display output for mime type(s): application/vnd.plotly.v1+json

```
optuna.visualization.plot_param_importances(study).show()
```

Unable to display output for mime type(s): application/vnd.plotly.v1+json

```
optuna.visualization.plot_parallel_coordinate(study).show()
```

Unable to display output for mime type(s): application/vnd.plotly.v1+json

### 9.3 Scaling Gradient Boosting for Large Datasets and Smarter Tuning

Check out `HistGradientBoostingRegressor()` and `HistGradientBoostingClassifier()` for a faster gradient boosting algorithm for big datasets (*more than 10,000 observations*).

Check out tips for faster hyperparameter tuning, such as tuning `max_leaf_nodes` instead of `max_depth` [here](#).



## 9.4 Independent Study

In this notebook, we used the car dataset for a guided regression task to illustrate the core hyperparameters in gradient boosting and how to tune them to balance bias and variance.

For your practice, please work with the **diabetes dataset** and complete the following:

- Fit a baseline gradient boosting classifier.
- Tune key hyperparameters: `learning_rate`, `n_estimators`, `max_depth`, and `subsample`.
- Use **early stopping** to determine the optimal number of trees.
- Compare training and test `roc_auc` before and after tuning.
- Visualize the learning curve (training vs test error across iterations).
- Summarize what combination of hyperparameters yielded the best performance and how they impacted bias and variance.

Feel free to use `GridSearchCV`, `BayesSearchCV`, or other tuning as you prefer.

```
train = pd.read_csv('./Datasets/diabetes_train.csv')
test = pd.read_csv('./Datasets/diabetes_test.csv')
```

```
print(train.shape, test.shape)
train.head()
```

(614, 9) (154, 9)

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	2	88	74	19	53	29.0	0.229	22
1	2	129	84	0	0	28.0	0.284	27
2	0	102	78	40	90	34.5	0.238	24
3	0	123	72	0	0	36.3	0.258	52
4	1	144	82	46	180	46.1	0.335	46

```
# check the distribution of the target variable
train['Outcome'].value_counts(normalize=True)
```

```
Outcome
0    0.662866
1    0.337134
Name: proportion, dtype: float64
```

# 10 XGBoost

XGBoost is a very recently developed algorithm (2016). Thus, it's not yet there in standard textbooks. Here are some resources for it.

[Documentation](#)

[Slides](#)

[Reference paper](#)

[Video by author \(Tianqi Chen\)](#)

[Video by StatQuest](#)

## 10.1 Hyperparameters

The following are some of the important hyperparameters to tune in XGBoost:

1. Number of trees (`n_estimators`)
2. Depth of each tree (`max_depth`)
3. Learning rate (`learning_rate`)
4. Sampling observations / predictors (`subsample` for observations, `colsample_bytree` for predictors)
5. Regularization parameters (`reg_lambda` & `gamma`)

However, there are other hyperparameters that can be tuned as well. Check out the list of all hyperparameters in the XGBoost [documentation](#).

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import cross_val_score, train_test_split, KFold, cross_val_predict
from sklearn.metrics import mean_squared_error, r2_score, roc_curve, auc, precision_recall_curve
```

```

recall_score, precision_score, confusion_matrix
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV, ParameterGrid, StratifiedKFold, RandomizedSearchCV
from sklearn.ensemble import VotingRegressor, VotingClassifier, StackingRegressor, StackingClassifier
from sklearn.linear_model import LinearRegression, LogisticRegression, LassoCV, RidgeCV, ElasticNetCV
from sklearn.neighbors import KNeighborsRegressor
import itertools as it
import time as time
import xgboost as xgb
from skopt import BayesSearchCV
from skopt.space import Real, Categorical, Integer
from skopt.plots import plot_objective, plot_histogram, plot_convergence
import warnings
from IPython import display

```

```

#Using the same datasets as used for linear regression in STAT303-2,
#so that we can compare the non-linear models with linear regression
trainf = pd.read_csv('./Datasets/Car_features_train.csv')
trainp = pd.read_csv('./Datasets/Car_prices_train.csv')
testf = pd.read_csv('./Datasets/Car_features_test.csv')
testp = pd.read_csv('./Datasets/Car_prices_test.csv')
train = pd.merge(trainf, trainp)
test = pd.merge(testf, testp)
train.head()

```

	carID	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	18473	bmw	6 Series	2020	Semi-Auto	11	Diesel	145	53.3282	3.0	37980
1	15064	bmw	6 Series	2019	Semi-Auto	10813	Diesel	145	53.0430	3.0	33980
2	18268	bmw	6 Series	2020	Semi-Auto	6	Diesel	145	53.4379	3.0	36850
3	18480	bmw	6 Series	2017	Semi-Auto	18895	Diesel	145	51.5140	3.0	25998
4	18492	bmw	6 Series	2015	Automatic	62953	Diesel	160	51.4903	3.0	18990

```

X = train[['mileage', 'mpg', 'year', 'engineSize']]
Xtest = test[['mileage', 'mpg', 'year', 'engineSize']]
y = train['price']
ytest = test['price']

```

## 10.2 XGBoost for regression

### 10.2.1 Number of trees vs cross validation error

As the number of trees increase, the prediction bias will decrease. Like gradient boosting is relatively robust (*as compared to AdaBoost*) to over-fitting (why?) so a large number usually results in better performance. Note that the number of trees still need to be tuned for optimal performance.

```
def get_models():
    models = dict()
    # define number of trees to consider
    n_trees = [5, 10, 50, 100, 500, 1000, 2000, 5000]
    for n in n_trees:
        models[str(n)] = xgb.XGBRegressor(n_estimators=n, random_state=1)
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=5, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = np.sqrt(-cross_val_score(model, X, y, scoring='neg_mean_squared_error', cv=cv, n_jobs=-1))
    return scores

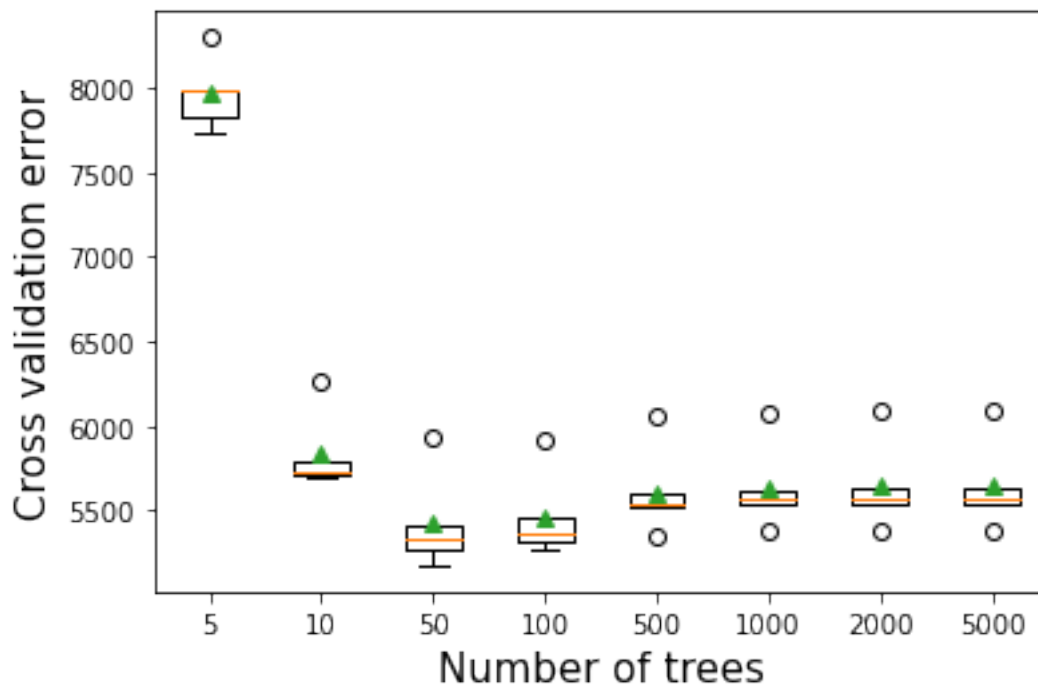
# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X, y)
    # store the results
    results.append(scores)
    names.append(name)
    # summarize the performance along the way
    print('>%s %.3f (%.3f)' % (name, np.mean(scores), np.std(scores)))
# plot model performance for comparison
plt.boxplot(results, labels=names, showmeans=True)
plt.ylabel('Cross validation error', fontsize=15)
plt.xlabel('Number of trees', fontsize=15)
```

```

>5 7961.485 (192.906)
>10 5837.134 (217.986)
>50 5424.788 (263.890)
>100 5465.396 (237.938)
>500 5608.350 (235.903)
>1000 5635.159 (236.664)
>2000 5642.669 (236.192)
>5000 5643.411 (236.074)

```

```
Text(0.5, 0, 'Number of trees')
```



### 10.2.2 Depth of tree vs cross validation error

As the depth of each weak learner (decision tree) increases, the complexity of the weak learner will increase. As the complexity increases, the prediction bias will decrease, while the prediction variance will increase. Thus, there will be an optimal depth of each weak learner that minimizes the prediction error.

```

# get a list of models to evaluate
def get_models():
    models = dict()
    # explore depths from 1 to 10
    for i in range(1,21):
        # define ensemble model
        models[str(i)] = xgb.XGBRegressor(random_state=1,max_depth=i)
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=10, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = np.sqrt(-cross_val_score(model, X, y, scoring='neg_mean_squared_error', cv=cv, n_jobs=-1))
    return scores

# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X, y)
    # store the results
    results.append(scores)
    names.append(name)
    # summarize the performance along the way
    print('>%s %.3f (%.3f)' % (name, np.mean(scores), np.std(scores)))
plt.boxplot(results, labels=names, showmeans=True)
plt.ylabel('Cross validation error',fontsize=15)
plt.xlabel('Depth of each tree',fontsize=15)

```

```

>1 7541.827 (545.951)
>2 6129.425 (393.357)
>3 5647.783 (454.318)
>4 5438.481 (453.726)
>5 5358.074 (379.431)
>6 5281.675 (383.848)
>7 5495.163 (459.356)
>8 5399.145 (380.437)
>9 5469.563 (384.004)

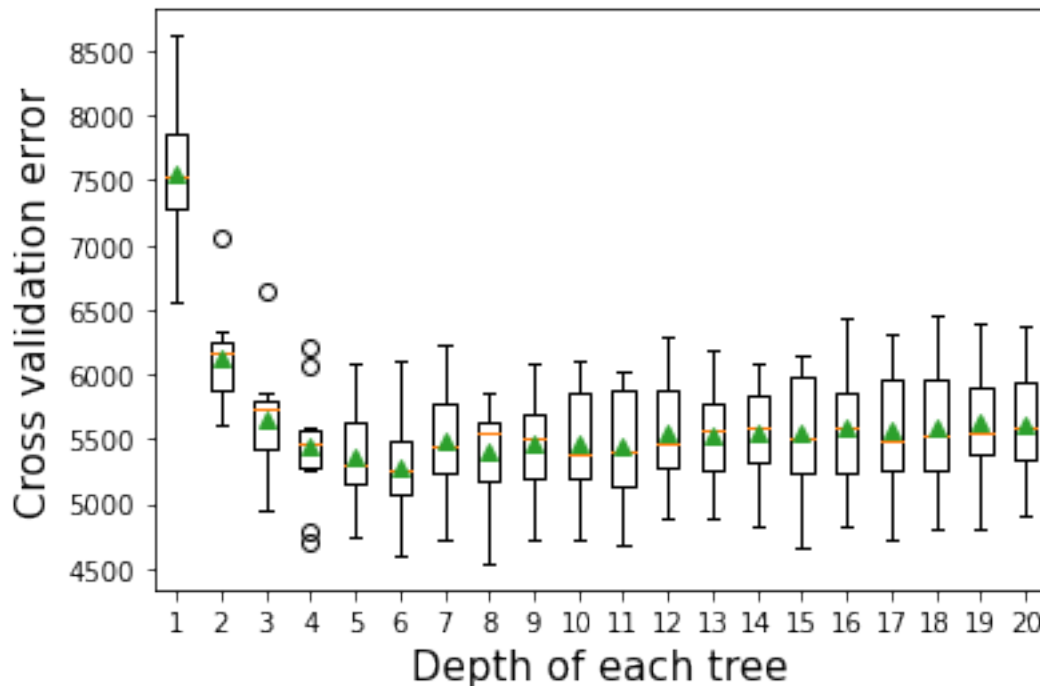
```

```

>10 5461.549 (416.630)
>11 5443.210 (432.863)
>12 5546.447 (412.097)
>13 5532.414 (369.131)
>14 5556.761 (362.746)
>15 5540.366 (452.612)
>16 5586.004 (451.199)
>17 5563.137 (464.344)
>18 5594.919 (480.221)
>19 5641.226 (451.713)
>20 5616.462 (417.405)

```

```
Text(0.5, 0, 'Depth of each tree')
```



### 10.2.3 Learning rate vs cross validation error

The optimal learning rate will depend on the number of trees, and vice-versa. If the learning rate is too low, it will take several trees to “learn” the response. If the learning rate is high, the response will be “learned” quickly (with fewer) trees. Learning too quickly will be prone to overfitting, while learning too slowly will be computationally expensive. Thus, there will be an optimal learning rate to minimize the prediction error.

```

def get_models():
    models = dict()
    # explore learning rates from 0.1 to 2 in 0.1 increments
    for i in [0.01,0.05,0.1,0.2,0.3,0.4,0.5,0.6,0.8,1.0]:
        key = '%.4f' % i
        models[key] = xgb.XGBRegressor(learning_rate=i,random_state=1)
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=10, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = np.sqrt(-cross_val_score(model, X, y, scoring='neg_mean_squared_error', cv=cv, n_jobs=-1))
    return scores

# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X, y)
    # store the results
    results.append(scores)
    names.append(name)
    # summarize the performance along the way
    print('>%s %.1f (%.1f)' % (name, np.mean(scores), np.std(scores)))
# plot model performance for comparison
plt.figure(figsize=(7, 7))
plt.boxplot(results, labels=names, showmeans=True)
plt.ylabel('Cross validation error',fontsize=15)
plt.xlabel('Learning rate',fontsize=15)

```

```

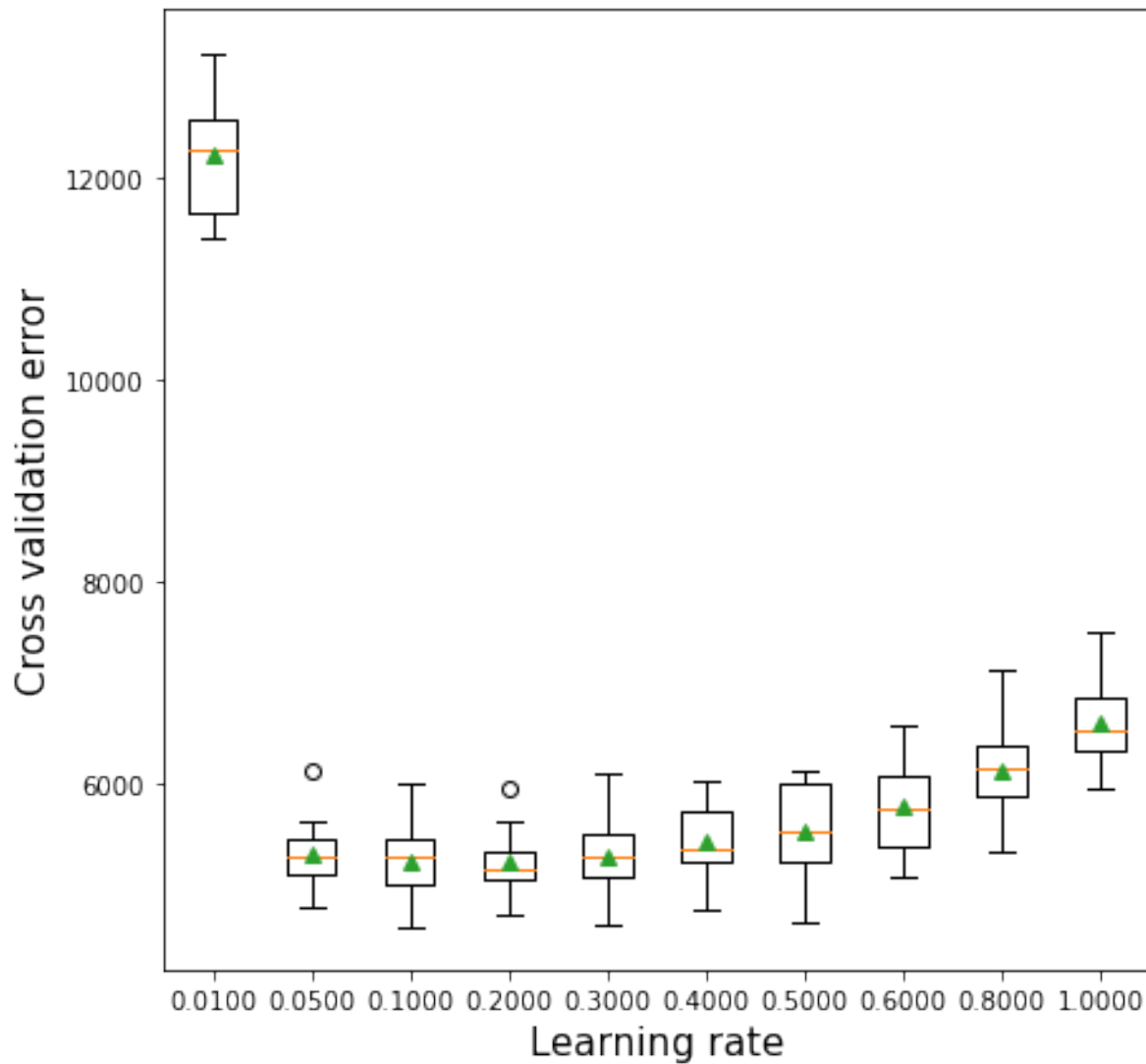
>0.0100 12223.8 (636.7)
>0.0500 5298.5 (383.5)
>0.1000 5236.3 (397.5)
>0.2000 5221.5 (347.5)
>0.3000 5281.7 (383.8)
>0.4000 5434.1 (364.6)
>0.5000 5537.0 (471.9)
>0.6000 5767.4 (478.5)

```



```
>0.8000 6132.7 (472.5)
>1.0000 6593.6 (408.9)
```

```
Text(0.5, 0, 'Learning rate')
```



#### 10.2.4 Regularization (reg\_lambda) vs cross validation error

The parameter `reg_lambda` penalizes the  $L2$  norm of the leaf scores. For example, in case of classification, it will penalize the summation of the square of log odds of the predicted

probability. This penalization will tend to reduce the log odds, thereby reducing the tendency to overfit. “*Reducing the log odds*” in layman terms will mean not being overly sure about the prediction.

Without regularization, the algorithm will be closer to the gradient boosting algorithm. Regularization may provide some additional boost to prediction accuracy by reducing over-fitting. In the example below, regularization with *reg\_lambda=1 turns out to be better than no regularization* (reg\_lambda=0)\*. Of course, too much regularization may increase bias so much such that it leads to a decrease in prediction accuracy.

```
def get_models():
    models = dict()
    # explore 'reg_lambda' from 0.1 to 2 in 0.1 increments
    for i in [0,0.5,1.0,1.5,2,10,100]:
        key = '%.4f' % i
        models[key] = xgb.XGBRegressor(reg_lambda=i,random_state=1)
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=10, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = np.sqrt(-cross_val_score(model, X, y, scoring='neg_mean_squared_error', cv=cv, n_jobs=-1))
    return scores

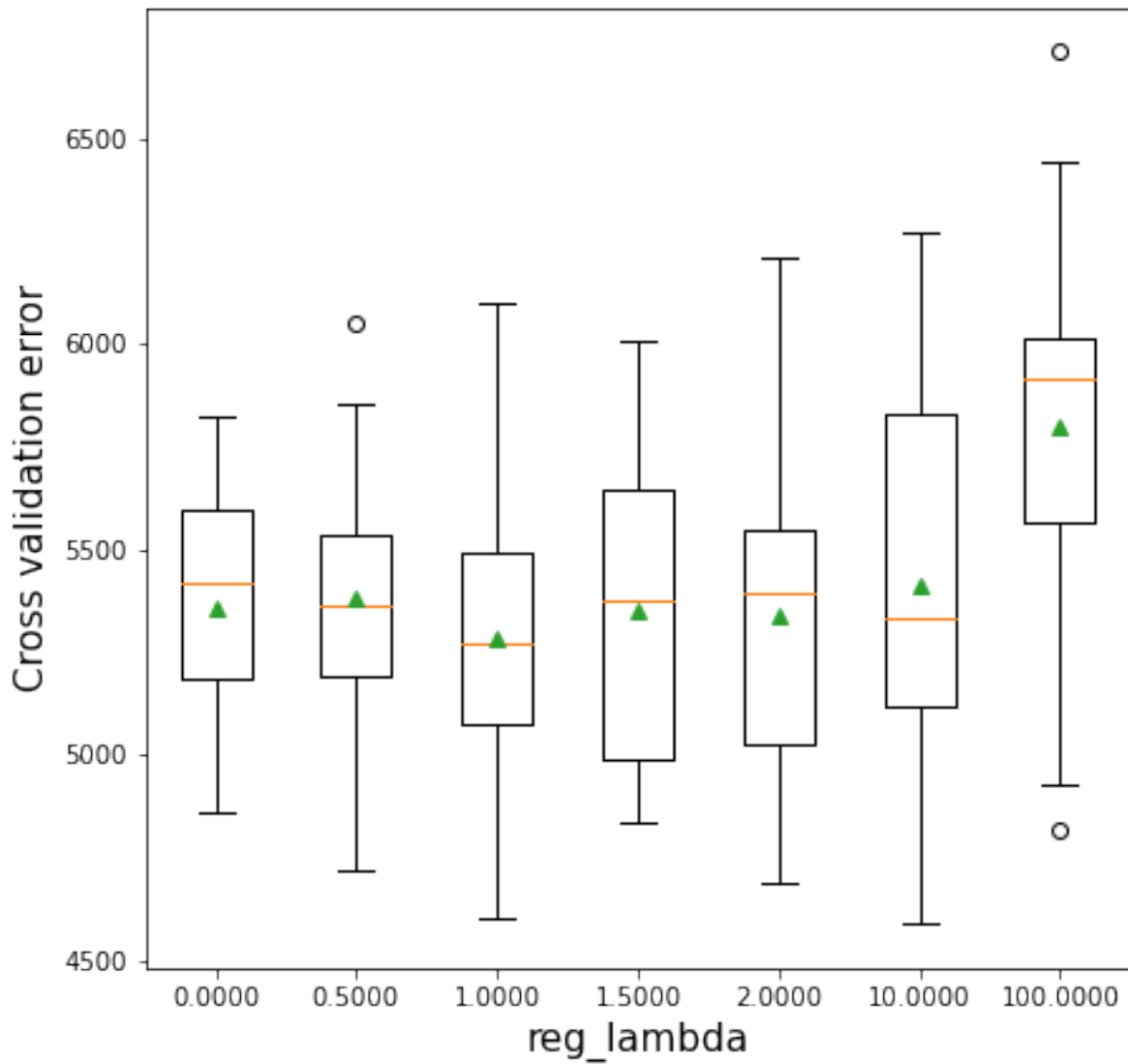
# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X, y)
    # store the results
    results.append(scores)
    names.append(name)
    # summarize the performance along the way
    print('>%s %.1f (%.1f)' % (name, np.mean(scores), np.std(scores)))
# plot model performance for comparison
plt.figure(figsize=(7, 7))
plt.boxplot(results, labels=names, showmeans=True)
plt.ylabel('Cross validation error',fontsize=15)
plt.xlabel('reg_lambda',fontsize=15)
```

```

>0.0000 5359.2 (317.0)
>0.5000 5382.7 (363.1)
>1.0000 5281.7 (383.8)
>1.5000 5348.0 (383.9)
>2.0000 5336.4 (426.6)
>10.0000 5410.9 (521.9)
>100.0000 5801.1 (563.7)

```

```
Text(0.5, 0, 'reg_lambda')
```



## 10.2.5 Regularization (gamma) vs cross validation error

The parameter `gamma` penalizes the tree based on the number of leaves. This is similar to the parameter `alpha` of cost complexity pruning. As `gamma` increases, more leaves will be pruned. Note that the previous parameter `reg_lambda` penalizes the leaf score, but does not prune the tree.

Without regularization, the algorithm will be closer to the gradient boosting algorithm. Regularization may provide some additional boost to prediction accuracy by reducing over-fitting. However, in the example below, no regularization (in terms of `gamma=0`) turns out to be better than a non-zero regularization. (*reg\_lambda=0*).

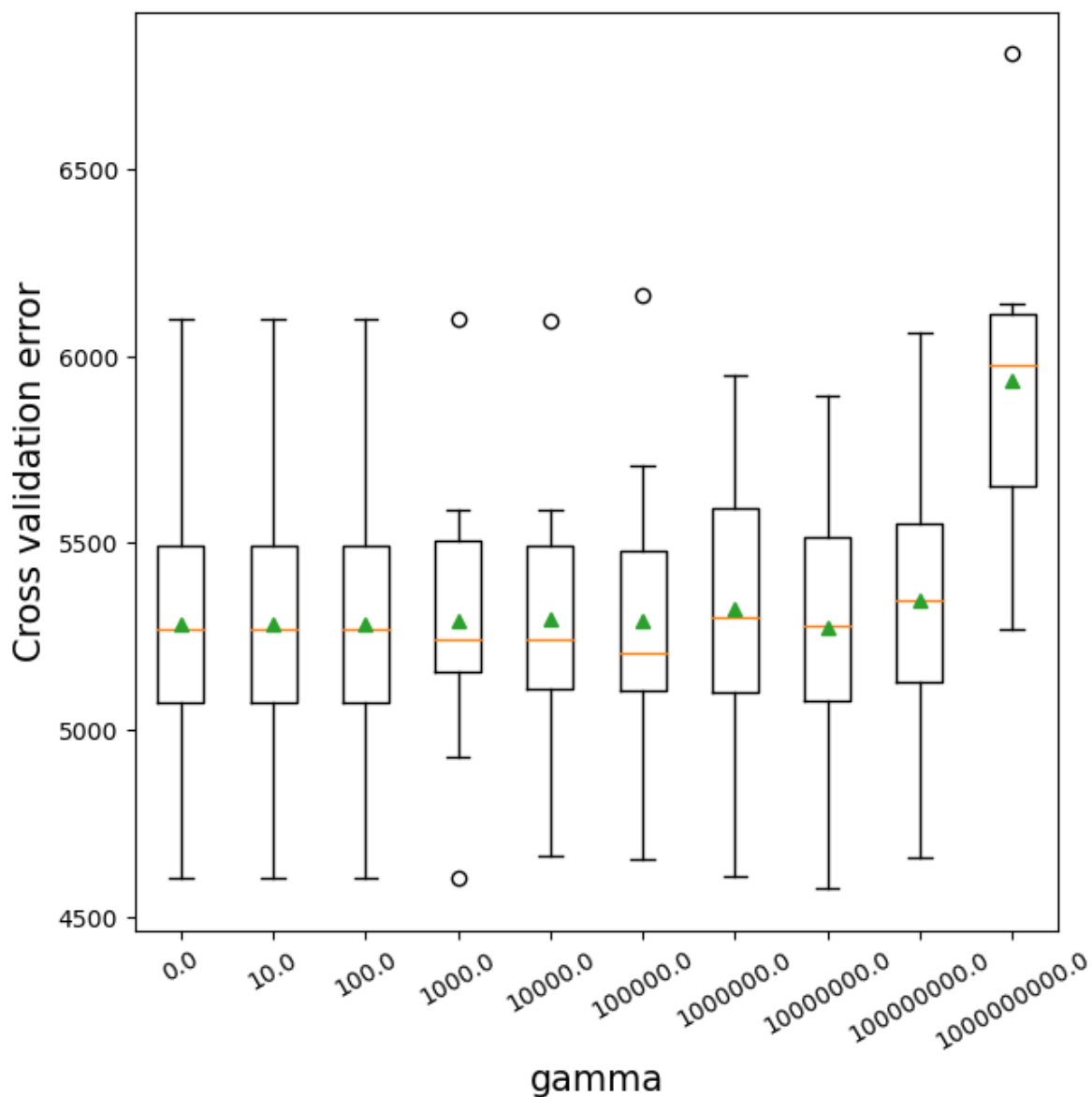
```
def get_models():
    models = dict()
    # explore gamma from 0.1 to 2 in 0.1 increments
    for i in [0,10,1e2,1e3,1e4,1e5,1e6,1e7,1e8,1e9]:
        key = '%.4f' % i
        models[key] = xgb.XGBRegressor(gamma=i,random_state=1)
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = KFold(n_splits=10, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores = np.sqrt(-cross_val_score(model, X, y, scoring='neg_mean_squared_error', cv=cv, n_jobs=-1))
    return scores

# get the models to evaluate
models = get_models()
# evaluate the models and store results
results, names = list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X, y)
    # store the results
    results.append(scores)
    names.append(name)
    # summarize the performance along the way
    print('>%s %.1f (%.1f)' % (name, np.mean(scores), np.std(scores)))
# plot model performance for comparison
plt.figure(figsize=(7, 7))
plt.boxplot(results, labels=names, showmeans=True)
```

```
plt.ylabel('Cross validation error',fontsize=15)
plt.xlabel('gamma',fontsize=15)
#ax.set_xticklabels(x.astype(int))
plt.xticks(ticks=plt.xticks()[0].astype(int), labels=np.round([0,10,1e2,1e3,1e4,1e5,1e6,1e7,
rotation = 30);
```

```
>0.0000 5281.7 (383.8)
>10.0000 5281.7 (383.8)
>100.0000 5281.7 (383.8)
>1000.0000 5291.8 (381.8)
>10000.0000 5295.7 (370.2)
>100000.0000 5293.0 (402.5)
>1000000.0000 5322.2 (368.9)
>10000000.0000 5273.7 (409.8)
>100000000.0000 5345.4 (373.9)
>1000000000.0000 5932.3 (397.6)
```



### 10.2.6 Tuning XGboost regressor

Along with `max_depth`, `learning_rate`, and `n_estimators`, here we tune `reg_lambda` - the regularization parameter for penalizing the tree predictions.

```
#K-fold cross validation to find optimal parameters for XGBoost
start_time = time.time()
param_grid = {'max_depth': [4,6,8],
```

```

        'learning_rate': [0.01, 0.05, 0.1],
        'reg_lambda': [0, 1, 10],
        'n_estimators': [100, 500, 1000],
        'gamma': [0, 10, 100],
        'subsample': [0.5, 0.75, 1.0],
        'colsample_bytree': [0.5, 0.75, 1.0]}

cv = KFold(n_splits=5, shuffle=True, random_state=1)
optimal_params = RandomizedSearchCV(estimator=xgb.XGBRegressor(random_state=1),
                                    param_distributions = param_grid, n_iter = 200,
                                    verbose = 1,
                                    n_jobs=-1,
                                    cv = cv)

optimal_params.fit(X,y)
print("Optimal parameter values =", optimal_params.best_params_)
print("Optimal cross validation R-squared = ", optimal_params.best_score_)
print("Time taken = ", round((time.time()-start_time)/60), " minutes")

```

Fitting 5 folds for each of 200 candidates, totalling 1000 fits

Optimal parameter values = {'subsample': 0.75, 'reg\_lambda': 1, 'n\_estimators': 1000, 'max\_d

Optimal cross validation R-squared = 0.9002580404500382

Time taken = 4 minutes

#RMSE based on the optimal parameter values

```
np.sqrt(mean_squared_error(optimal_params.best_estimator_.predict(Xtest), ytest))
```

5497.553788113875

Let us use Bayes search to tune the model.

```

model = xgb.XGBRegressor(random_state = 1)

grid = {'max_leaves': Integer(4, 5000),
        'learning_rate': Real(0.0001, 1.0),
        'reg_lambda': Real(0, 1e4),
        'n_estimators': Integer(2, 2000),
        'gamma': Real(0, 1e11),
        'subsample': Real(0.1, 1.0),
        'colsample_bytree': Real(0.1, 1.0)}

```

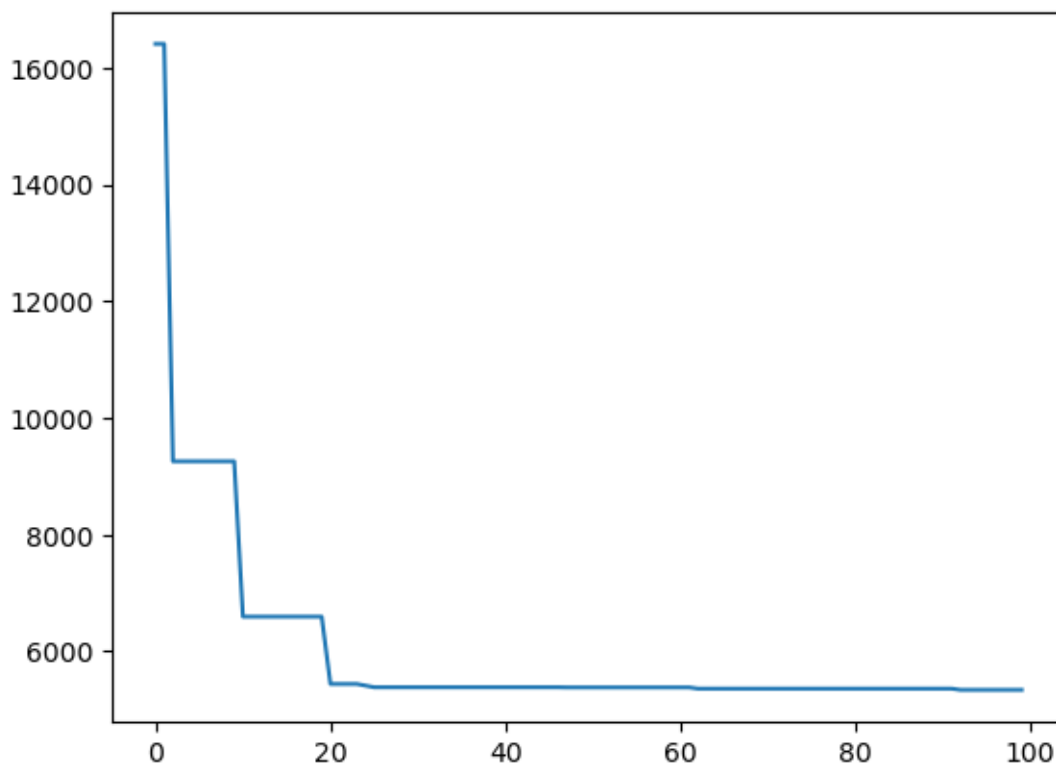
```

kfold = KFold(n_splits = 5, shuffle = True, random_state = 1)
gcv = BayesSearchCV(model, search_spaces = grid, cv = kfold, n_iter = 100, random_state = 1,
                    scoring = 'neg_root_mean_squared_error', n_jobs = -1)
paras = list(gcv.search_spaces.keys())
paras.sort()

def monitor(optim_result):
    cv_values = pd.Series(optim_result['func_vals']).cummin()
    display.clear_output(wait = True)
    min_ind = pd.Series(optim_result['func_vals']).argmin()
    print(paras, "=", optim_result['x_iters'][min_ind], pd.Series(optim_result['func_vals']))
    sns.lineplot(cv_values)
    plt.show()
gcv.fit(X, y, callback = monitor)

```

['colsample\_bytree', 'gamma', 'learning\_rate', 'max\_leaves', 'n\_estimators', 'reg\_lambda', '...



BayesSearchCV(cv=KFold(n\_splits=5, random\_state=1, shuffle=True),



```

estimator=XGBRegressor(base_score=None, booster=None,
                        callbacks=None, colsample_bylevel=None,
                        colsample_bynode=None,
                        colsample_bytree=None,
                        early_stopping_rounds=None,
                        enable_categorical=False, eval_metric=None,
                        feature_types=None, gamma=None,
                        gpu_id=None, grow_policy=None,
                        importance_type=None,
                        inte...
                        'learning_rate': Real(low=0.0001, high=1.0, prior='uniform', tra
                        'max_leaves': Integer(low=4, high=5000, prior='uniform', trans
                        'n_estimators': Integer(low=2, high=2000, prior='uniform', tran
                        'reg_lambda': Real(low=0, high=10000.0, prior='uniform', trans
                        'subsample': Real(low=0.1, high=1.0, prior='uniform', transform

```

```

model1 = xgb.XGBRegressor(random_state = 1, colsample_bytree = 0.85, gamma = 0, learning_rate = 0.1,
                           max_leaves = 802, n_estimators = 1023, reg_lambda = 1394, subsample = 0.8)

```

```

np.sqrt(mean_squared_error(model1.predict(Xtest),ytest))

```

5466.076861800755

We got a different set of optimal hyperparameters with Bayes search. Thus, ensembling the model based on the two sets of hyperparameters is likely to improve the accuracy over the individual models.

```

model2 = xgb.XGBRegressor(random_state = 1, colsample_bytree = 1.0, gamma = 100, learning_rate = 0.1,
                           max_depth = 8, n_estimators = 1000, reg_lambda = 1, subsample = 0.8)

```

```

np.sqrt(mean_squared_error(0.5*model1.predict(Xtest)+0.5*model2.predict(Xtest),ytest))

```

5393.379834226845

## 10.2.7 Early stopping with XGBoost

If we have a test dataset (*or we can further split the train data into a smaller train and test data*), we can use it with the `early_stopping_rounds` argument of XGBoost, where it will stop growing trees once the model accuracy fails to increase for a certain number of consecutive iterations, given as `early_stopping_rounds`.

```
X_train_sub, X_test_sub, y_train_sub, y_test_sub = \
train_test_split(X, y, test_size = 0.2, random_state = 45)
```

```
model = xgb.XGBRegressor(random_state = 1, max_depth = 8, learning_rate = 0.01,
                        n_estimators = 20000, reg_lambda = 1, gamma = 100, subsample = 0.75,
model.fit(X_train_sub, y_train_sub, eval_set = ((X_test_sub, y_test_sub))), early_stopping_
```

The results of the code are truncated to save space. A snapshot of the beginning and end of the results is below. The algorithm keeps adding trees to the model until the RMSE ceases to decrease for 250 consecutive iterations.

```
<IPython.core.display.Image object>
```

```
print("XGBoost RMSE = ", np.sqrt(mean_squared_error(model.predict(Xtest), ytest)))
```

```
XGBoost RMSE = 5508.787454011525
```

Let us further reduce the learning rate to 0.001 and see if the accuracy increases further on the test data. We'll use the `early_stopping_rounds` argument to stop growing trees once the accuracy fails to increase for 250 consecutive iterations.

```
model = xgb.XGBRegressor(random_state = 1, max_depth = 8, learning_rate = 0.001,
                        n_estimators = 20000, reg_lambda = 1, gamma = 100, subsample = 0.75,
model.fit(X_train_sub, y_train_sub, eval_set = ((X_test_sub, y_test_sub))), early_stopping_
```

```
<IPython.core.display.Image object>
```

```
print("XGBoost RMSE = ", np.sqrt(mean_squared_error(model.predict(Xtest), ytest)))
```

```
XGBoost RMSE = 5483.518711988693
```

Note that the accuracy on this test data has further increased with a lower learning rate.

Let us combine the XGBoost model with other tuned models from earlier chapters.

```

#Tuned AdaBoost model from Section 7.2.4
model_ada = AdaBoostRegressor(base_estimator=DecisionTreeRegressor(max_depth=10),n_estimators=100,
                              random_state=1).fit(X,y)
print("AdaBoost RMSE = ", np.sqrt(mean_squared_error(model_ada.predict(Xtest),ytest)))

#Tuned Random forest model from Section 6.1.2
model_rf = RandomForestRegressor(n_estimators=300, random_state=1,
                                n_jobs=-1, max_features=2).fit(X, y)
print("Random Forest RMSE = ",np.sqrt(mean_squared_error(model_rf.predict(Xtest),ytest)))

#Tuned gradient boosting model from Section 8.2.5
model_gb = GradientBoostingRegressor(max_depth=8,n_estimators=100,learning_rate=0.1,
                                     random_state=1,loss='huber').fit(X,y)
print("Gradient boost RMSE = ",np.sqrt(mean_squared_error(model_gb.predict(Xtest),ytest)))

```

```

AdaBoost RMSE = 5693.165811600585
Random Forest RMSE = 5642.45839697972
Gradient boost RMSE = 5405.787029062213

```

```

#Ensemble model
pred_xgb = model.predict(Xtest)    #XGBoost
pred_ada = model_ada.predict(Xtest)#AdaBoost
pred_rf = model_rf.predict(Xtest)  #Random Forest
pred_gb = model_gb.predict(Xtest)  #Gradient boost
pred = 0.25*pred_xgb + 0.25*pred_ada + 0.25*pred_rf + 0.25*pred_gb #Option 1 - All models are given equal weight
#pred = 0.15*pred1+0.15*pred2+0.15*pred3+0.55*pred4 #Option 2 - Higher weight to the better model
print("Ensemble model RMSE = ", np.sqrt(mean_squared_error(pred,ytest)))

```

```

Ensemble model RMSE = 5352.145010078119

```

Combined, the random forest model, gradient boost, XGBoost and the Adaboost model do better than each of the individual models.

## 10.3 XGBoost for classification

```

data = pd.read_csv('./Datasets/Heart.csv')
data.dropna(inplace = True)
data.head()

```

	Age	Sex	ChestPain	RestBP	Chol	Fbs	RestECG	MaxHR	ExAng	Oldpeak	Slope	Ca
0	63	1	typical	145	233	1	2	150	0	2.3	3	0.0
1	67	1	asymptomatic	160	286	0	2	108	1	1.5	2	3.0
2	67	1	asymptomatic	120	229	0	2	129	1	2.6	2	2.0
3	37	1	nonanginal	130	250	0	0	187	0	3.5	3	0.0
4	41	0	nontypical	130	204	0	2	172	0	1.4	1	0.0

```
#Response variable
y = pd.get_dummies(data['AHD'])['Yes']

#Creating a dataframe for predictors with dummy variables replacing the categorical variables
X = data.drop(columns = ['AHD','ChestPain','Thal'])
X = pd.concat([X,pd.get_dummies(data['ChestPain']),pd.get_dummies(data['Thal'])],axis=1)
X.head()
```

	Age	Sex	RestBP	Chol	Fbs	RestECG	MaxHR	ExAng	Oldpeak	Slope	Ca	asymptomatic
0	63	1	145	233	1	2	150	0	2.3	3	0.0	0
1	67	1	160	286	0	2	108	1	1.5	2	3.0	1
2	67	1	120	229	0	2	129	1	2.6	2	2.0	1
3	37	1	130	250	0	0	187	0	3.5	3	0.0	0
4	41	0	130	204	0	2	172	0	1.4	1	0.0	0

```
#Creating train and test datasets
Xtrain,Xtest,ytrain,ytest = train_test_split(X,y,train_size = 0.5,random_state=1)
```

XGBoost has an additional parameter for classification: **scale\_pos\_weight**

Gradients are used as the basis for fitting subsequent trees added to boost or correct errors made by the existing state of the ensemble of decision trees.

The **scale\_pos\_weight** value is used to scale the gradient for the positive class.

This has the effect of scaling errors made by the model during training on the positive class and encourages the model to over-correct them. In turn, this can help the model achieve better performance when making predictions on the positive class. Pushed too far, it may result in the model overfitting the positive class at the cost of worse performance on the negative class or both classes.

As such, the **scale\_pos\_weight** hyperparameter can be used to train a class-weighted or cost-sensitive version of XGBoost for imbalanced classification.

A sensible default value to set for the `scale_pos_weight` hyperparameter is the inverse of the class distribution. For example, for a dataset with a 1 to 100 ratio for examples in the minority to majority classes, the `scale_pos_weight` can be set to 100. This will give classification errors made by the model on the minority class (positive class) 100 times more impact, and in turn, 100 times more correction than errors made on the majority class.

## Reference

```
start_time = time.time()
param_grid = {'n_estimators': [25, 100, 500],
              'max_depth': [6, 7, 8],
              'learning_rate': [0.01, 0.1, 0.2],
              'gamma': [0.1, 0.25, 0.5],
              'reg_lambda': [0, 0.01, 0.001],
              'scale_pos_weight': [1.25, 1.5, 1.75] #Control the balance of positive and negative
            }

cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=1)
optimal_params = GridSearchCV(estimator=xgb.XGBClassifier(objective = 'binary:logistic', random_state=1,
                                                         use_label_encoder=False),
                             param_grid = param_grid,
                             scoring = 'accuracy',
                             verbose = 1,
                             n_jobs=-1,
                             cv = cv)

optimal_params.fit(Xtrain, ytrain)
print(optimal_params.best_params_, optimal_params.best_score_)
print("Time taken = ", (time.time()-start_time)/60, " minutes")
```

Fitting 5 folds for each of 729 candidates, totalling 3645 fits

[22:00:02] WARNING: D:\bld\xgboost-split\_1645118015404\work\src\learner.cc:1115: Starting in {'gamma': 0.25, 'learning\_rate': 0.2, 'max\_depth': 6, 'n\_estimators': 25, 'reg\_lambda': 0.01

```
cv_results=pd.DataFrame(optimal_params.cv_results_)
cv_results.sort_values(by = 'mean_test_score', ascending=False)[0:5]
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_gamma	param_learning_rate
409	0.111135	0.017064	0.005629	0.000737	0.25	0.2
226	0.215781	0.007873	0.005534	0.001615	0.1	0.2
290	1.391273	0.107808	0.007723	0.006286	0.25	0.01
266	1.247463	0.053597	0.006830	0.002728	0.25	0.01

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_gamma	param_learning_rate
269	1.394361	0.087307	0.005530	0.001718	0.25	0.01

```
#Function to compute confusion matrix and prediction accuracy on test/train data
def confusion_matrix_data(data,actual_values,model,cutoff=0.5):
#Predict the values using the Logit model
    pred_values = model.predict_proba(data)[:,-1]
# Specify the bins
    bins=np.array([0,cutoff,1])
#Confusion matrix
    cm = np.histogram2d(actual_values, pred_values, bins=bins)[0]
    cm_df = pd.DataFrame(cm)
    cm_df.columns = ['Predicted 0','Predicted 1']
    cm_df = cm_df.rename(index={0: 'Actual 0',1:'Actual 1'})
# Calculate the accuracy
    accuracy = 100*(cm[0,0]+cm[1,1])/cm.sum()
    fnr = 100*(cm[1,0])/(cm[1,0]+cm[1,1])
    precision = 100*(cm[1,1])/(cm[0,1]+cm[1,1])
    fpr = 100*(cm[0,1])/(cm[0,0]+cm[0,1])
    tpr = 100*(cm[1,1])/(cm[1,0]+cm[1,1])
    print("Accuracy = ", accuracy)
    print("Precision = ", precision)
    print("FNR = ", fnr)
    print("FPR = ", fpr)
    print("TPR or Recall = ", tpr)
    print("Confusion matrix = \n", cm_df)
    return (" ")
```

```
model4 = xgb.XGBClassifier(objective = 'binary:logistic',random_state=1,gamma=0.25,learning_rate=0.01,
                           n_estimators = 500,reg_lambda = 0.01,scale_pos_weight=1.75)
model4.fit(Xtrain,ytrain)
model4.score(Xtest,ytest)
```

0.7718120805369127

```
#Computing the accuracy
y_pred = model4.predict(Xtest)
print("Accuracy: ",accuracy_score(y_pred, ytest)*100)

#Computing the ROC-AUC
```

```

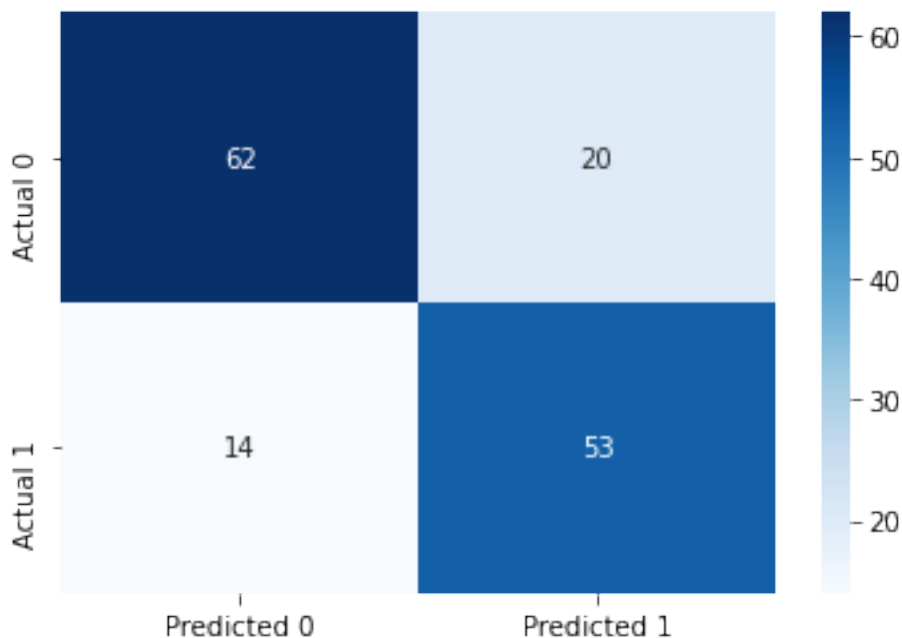
y_pred_prob = model4.predict_proba(Xtest)[:,-1]
fpr, tpr, auc_thresholds = roc_curve(ytest, y_pred_prob)
print("ROC-AUC: ",auc(fpr, tpr))# AUC of ROC

#Computing the precision and recall
print("Precision: ", precision_score(ytest, y_pred))
print("Recall: ", recall_score(ytest, y_pred))

#Confusion matrix
cm = pd.DataFrame(confusion_matrix(ytest, y_pred), columns=['Predicted 0', 'Predicted 1'],
                  index = ['Actual 0', 'Actual 1'])
sns.heatmap(cm, annot=True, cmap='Blues', fmt='g');

```

Accuracy: 77.18120805369128  
 ROC-AUC: 0.8815070986530761  
 Precision: 0.726027397260274  
 Recall: 0.7910447761194029



If we increase the value of `scale_pos_weight`, the model will focus on classifying positives more correctly. This will increase the recall (true positive rate) since the focus is on identifying all positives. However, this will lead to identifying positives aggressively, and observations 'similar' to observations of the positive class will also be predicted as positive resulting in an

increase in false positives and a decrease in precision. See the trend below as we increase the value of `scale_pos_weight`.

### 10.3.1 Precision & recall vs `scale_pos_weight`

```
def get_models():
    models = dict()
    # explore 'scale_pos_weight' from 0.1 to 2 in 0.1 increments
    for i in [0,1,10,1e2,1e3,1e4,1e5,1e6,1e7,1e8,1e9]:
        key = '%.0f' % i
        models[key] = xgb.XGBClassifier(objective = 'binary:logistic',scale_pos_weight=i,random_state=1)
    return models

# evaluate a given model using cross-validation
def evaluate_model(model, X, y):
    # define the evaluation procedure
    cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=1)
    # evaluate the model and collect the results
    scores_recall = cross_val_score(model, X, y, scoring='recall', cv=cv, n_jobs=-1)
    scores_precision = cross_val_score(model, X, y, scoring='precision', cv=cv, n_jobs=-1)
    return list([scores_recall,scores_precision])

# get the models to evaluate
models = get_models()
# evaluate the models and store results
results_recall, results_precision, names = list(), list(), list()
for name, model in models.items():
    # evaluate the model
    scores = evaluate_model(model, X, y)
    scores_recall = scores[0]
    scores_precision = scores[1]
    # store the results
    results_recall.append(scores_recall)
    results_precision.append(scores_precision)
    names.append(name)
    # summarize the performance along the way
    print('>%s %.2f (%.2f)' % (name, np.mean(scores_recall), np.std(scores_recall)))
# plot model performance for comparison
plt.figure(figsize=(7, 7))
sns.set(font_scale = 1.5)
pdata = pd.DataFrame(results_precision)
```



```

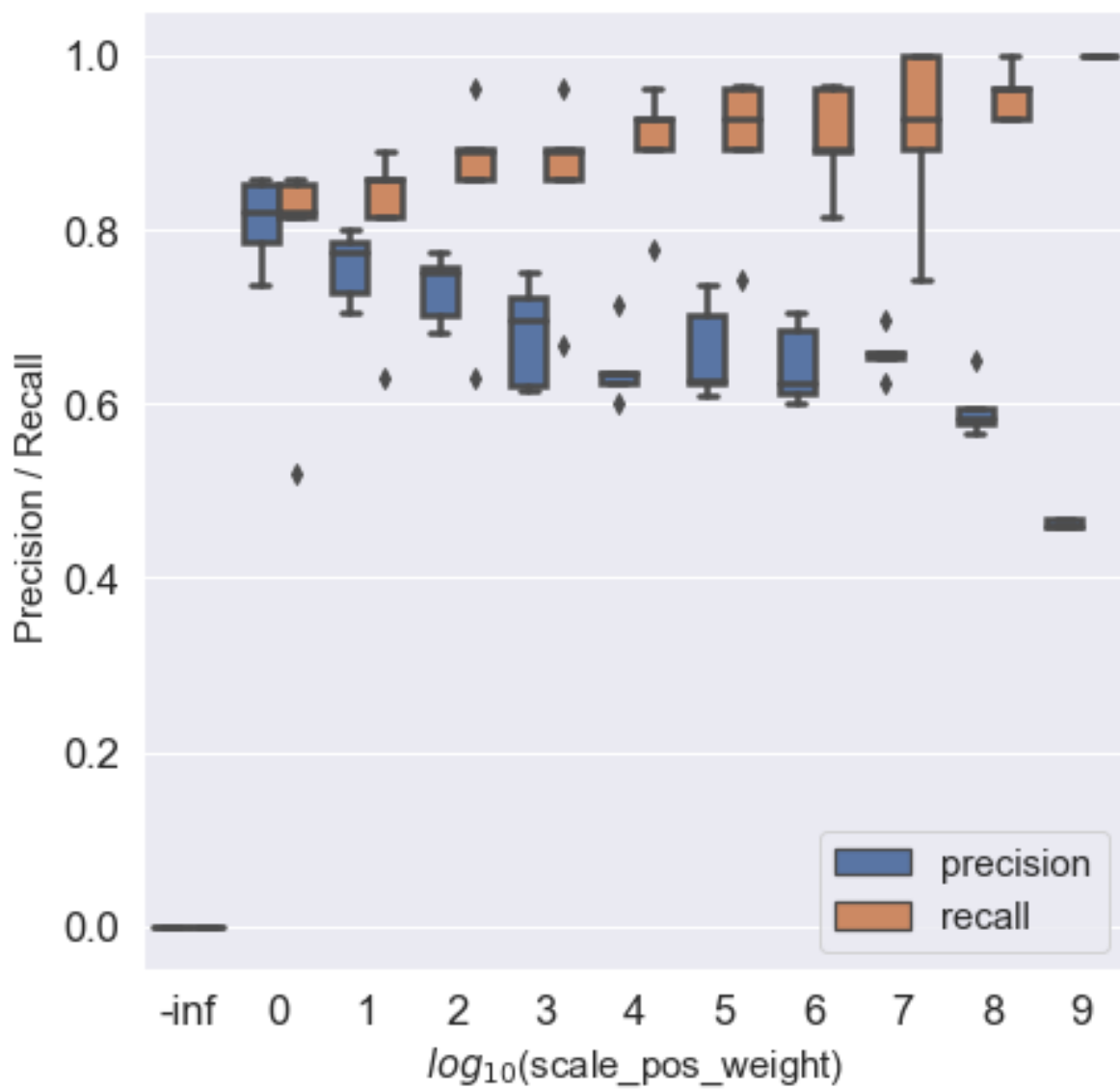
pdata.columns = list(['p1','p2','p3','p4','p5'])
pdata['metric'] = 'precision'
rdata = pd.DataFrame(results_recall)
rdata.columns = list(['p1','p2','p3','p4','p5'])
rdata['metric'] = 'recall'
pr_data = pd.concat([pdata,rdata])
pr_data.reset_index(drop=False,inplace= True)
#sns.boxplot(x="day", y="total_bill", hue="time",pr_data=tips, linewidth=2.5)
pr_data_melt=pr_data.melt(id_vars = ['index','metric'])
pr_data_melt['index']=pr_data_melt['index']-1
pr_data_melt['index'] = pr_data_melt['index'].astype('str')
pr_data_melt.replace(to_replace='-1',value =  '-inf',inplace=True)
sns.boxplot(x='index', y="value", hue="metric", data=pr_data_melt, linewidth=2.5)
plt.xlabel('$\log_{10}$(scale_pos_weight)',fontsize=15)
plt.ylabel('Precision / Recall ',fontsize=15)
plt.legend(loc="lower right", frameon=True, fontsize=15)

```

```

>0 0.00 (0.00)
>1 0.77 (0.13)
>10 0.81 (0.09)
>100 0.85 (0.11)
>1000 0.85 (0.10)
>10000 0.90 (0.06)
>100000 0.90 (0.08)
>1000000 0.90 (0.06)
>10000000 0.91 (0.10)
>100000000 0.96 (0.03)
>1000000000 1.00 (0.00)

```



# 11 LightGBM and CatBoost

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import cross_val_score, train_test_split, KFold, cross_val_predict
from sklearn.metrics import mean_squared_error, r2_score, roc_curve, auc, precision_recall_curve
recall_score, precision_score, confusion_matrix
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV, ParameterGrid, StratifiedKFold, RandomizedSearchCV
from sklearn.ensemble import VotingRegressor, VotingClassifier, StackingRegressor, StackingClassifier
from sklearn.linear_model import LinearRegression, LogisticRegression, LassoCV, RidgeCV, ElasticNetCV
from sklearn.neighbors import KNeighborsRegressor
import itertools as it
import time as time
import xgboost as xgb
from lightgbm import LGBMRegressor
from catboost import CatBoostRegressor
from skopt import BayesSearchCV
from skopt.space import Real, Categorical, Integer
from skopt.plots import plot_objective, plot_histogram, plot_convergence
import warnings
from IPython import display
```

We'll continue to use the same datasets that we have been using throughout the course.

```
#Using the same datasets as used for linear regression in STAT303-2,
#so that we can compare the non-linear models with linear regression
trainf = pd.read_csv('./Datasets/Car_features_train.csv')
trainp = pd.read_csv('./Datasets/Car_prices_train.csv')
testf = pd.read_csv('./Datasets/Car_features_test.csv')
testp = pd.read_csv('./Datasets/Car_prices_test.csv')
train = pd.merge(trainf, trainp)
```

```
test = pd.merge(testf, testp)
train.head()
```

	carID	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	18473	bmw	6 Series	2020	Semi-Auto	11	Diesel	145	53.3282	3.0	37980
1	15064	bmw	6 Series	2019	Semi-Auto	10813	Diesel	145	53.0430	3.0	33980
2	18268	bmw	6 Series	2020	Semi-Auto	6	Diesel	145	53.4379	3.0	36850
3	18480	bmw	6 Series	2017	Semi-Auto	18895	Diesel	145	51.5140	3.0	25998
4	18492	bmw	6 Series	2015	Automatic	62953	Diesel	160	51.4903	3.0	18990

```
X = train[['mileage', 'mpg', 'year', 'engineSize']]
Xtest = test[['mileage', 'mpg', 'year', 'engineSize']]
y = train['price']
ytest = test['price']
```

## 11.1 LightGBM

LightGBM is a gradient boosting decision tree algorithm developed by Microsoft in 2017. LightGBM outperforms XGBoost in terms of computational speed, and provides comparable accuracy in general. The following two key features in LightGBM that make it faster than XGBoost:

**1. Gradient-based One-Side Sampling (GOSS):** Recall, in gradient boosting, we fit trees on the gradient of the loss function (*refer the gradient boosting algorithm in section 10.10.2 of the book, [Elements of Statistical Learning](#)*):

$$r_m = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}.$$

Observations that correspond to relatively larger gradients contribute more to minimizing the loss function as compared to observations with smaller gradients. The algorithm down-samples the observations with small gradients, while selecting all the observations with large gradients. As observations with large gradients contribute the most to the reduction in loss function when considering a split, the accuracy of loss reduction estimate is maintained even with a reduced sample size. This leads to similar performance in terms of prediction accuracy while reducing computation speed due to reduction in sample size to fit trees.

**2. Exclusive feature bundling (EFB):** This is useful when there are a lot of predictors, but the predictor space is sparse, i.e., most of the values are zero for several predictors, and

the predictors rarely take non-zero values simultaneously. This can typically happen in case of a lot of dummy variables in the data. In such a case, the predictors are bundled to create a single predictor.

In the example below you can see that feature1 and feature2 are mutually exclusive. In order to achieve non overlapping buckets we add bundle size of feature1 to feature2. This makes sure that non zero data points of bundled features (feature1 and feature2) reside in different buckets. In feature\_bundle buckets 1 to 4 contains non zero instances of feature1 and buckets 5,6 contain non zero instances of feature2 ([Reference](#)).

feature1	feature2	feature_bundle
0	2	6
0	1	5
0	2	6
1	0	1
2	0	2
3	0	3
4	0	4

Read the [LightGBM paper](#) for more details.

### 11.1.1 LightGBM for regression

Let us tune a lightGBM model for regression for our problem of predicting car price. We'll use the function [LGBMRegressor](#). For classification problems, [LGBMClassifier](#) can be used. Note that we are using the GOSS algorithm to downsample observations with smaller gradients.

```
#K-fold cross validation to find optimal parameters for LightGBM regressor
start_time = time.time()
param_grid = {'num_leaves': [20, 31, 40],
              'learning_rate': [0.01, 0.05, 0.1],
              'reg_lambda': [0, 10, 100],
              'n_estimators': [100, 500, 1000],
              'reg_alpha': [0, 10, 100],
              'subsample': [0.5, 0.75, 1.0],
              'colsample_bytree': [0.5, 0.75, 1.0]}

cv = KFold(n_splits=5, shuffle=True, random_state=1)
optimal_params = RandomizedSearchCV(estimator=LGBMRegressor(boosting_type = 'goss'),
                                   param_distributions = param_grid, n_iter = 200,
                                   verbose = 1, scoring='neg_root_mean_squared_error',
```

```

n_jobs=-1,random_state=1,
cv = cv)

optimal_params.fit(X,y)
print("Optimal parameter values =", optimal_params.best_params_)
print("Optimal cross validation RMSE = ",optimal_params.best_score_)
print("Time taken = ", round((time.time()-start_time)/60), " minutes")

```

Fitting 5 folds for each of 200 candidates, totalling 1000 fits

Optimal parameter values = {'subsample': 1.0, 'reg\_lambda': 10, 'reg\_alpha': 0, 'num\_leaves':

Optimal cross validation R-squared = -5670.309021679375

Time taken = 1 minutes

```

#RMSE based on the optimal parameter values of a LighGBM Regressor model
np.sqrt(mean_squared_error(optimal_params.best_estimator_.predict(Xtest),ytest))

```

5614.374498193448

Note that downsampling of small-gradient observations leads to faster execution time, but potentially by compromising some accuracy. We can expect to improve the accuracy by increasing the `top_rate` or the `other_rate` hyperparameters, but at an increased computational cost. In the cross-validation below, we have increased the `top_rate` to 0.5 from the default value of 0.2.

```

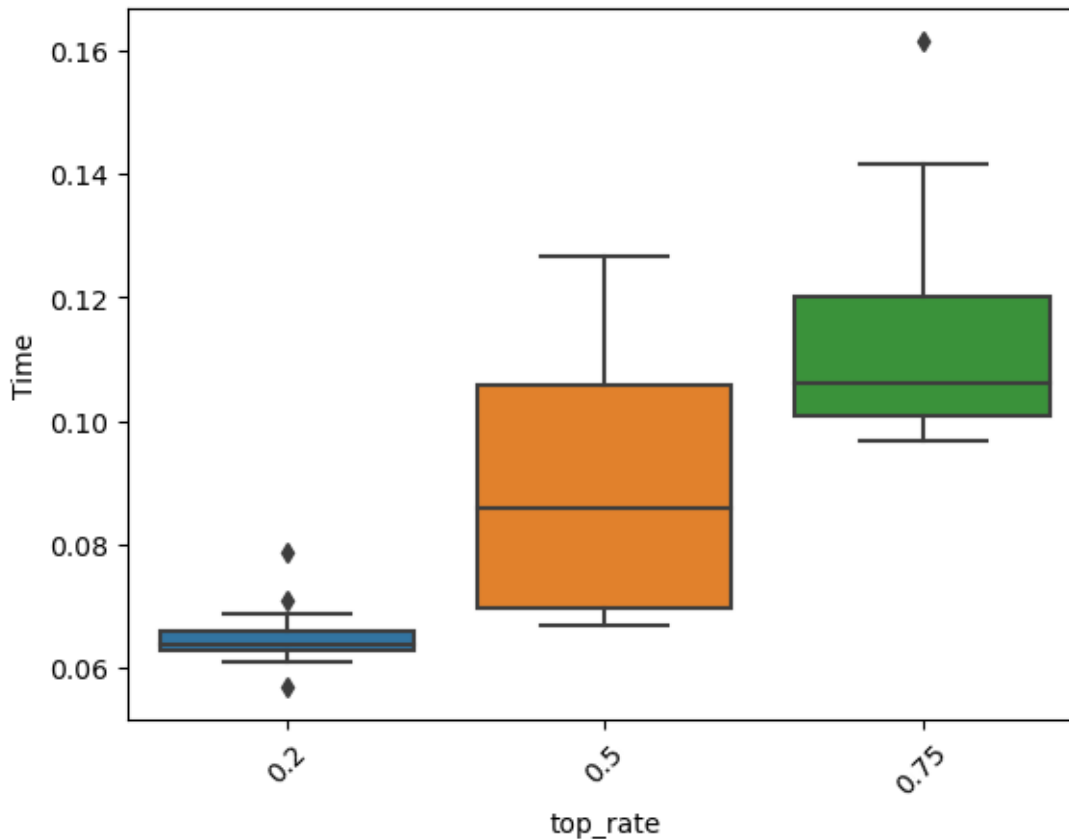
#K-fold cross validation to find optimal parameters for LightGBM regressor
start_time = time.time()
param_grid = {'num_leaves': [20, 31, 40],
              'learning_rate': [0.01, 0.05, 0.1],
              'reg_lambda': [0, 10, 100],
              'n_estimators': [100, 500, 1000],
              'reg_alpha': [0, 10, 100],
              'subsample': [0.5, 0.75, 1.0],
              'colsample_bytree': [0.5, 0.75, 1.0]}

cv = KFold(n_splits=5,shuffle=True,random_state=1)
optimal_params = RandomizedSearchCV(estimator=LGBMRegressor(boosting_type = 'goss', top_rate
                                param_distributions = param_grid, n_iter = 200,
                                verbose = 1, scoring='neg_root_mean_squared_error',
                                n_jobs=-1,random_state=1,
                                cv = cv)

optimal_params.fit(X,y)

```





### 11.1.2 LightGBM vs XGBoost

LightGBM model took 2 minutes for a random search with 1000 fits as compared to 7 minutes for an XGBoost model with 1000 fits on the same data (as shown below). In terms of prediction accuracy, we observe that the accuracy of LightGBM on test (*unseen*) data is comparable to that of XGBoost.

```
#K-fold cross validation to find optimal parameters for XGBoost
start_time = time.time()
param_grid = {'max_depth': [4,6,8],
              'learning_rate': [0.01, 0.05, 0.1],
              'reg_lambda': [0, 1, 10],
              'n_estimators': [100, 500, 1000],
              'gamma': [0, 10, 100],
              'subsample': [0.5, 0.75, 1.0],
              'colsample_bytree': [0.5, 0.75, 1.0]}
```



```

cv = KFold(n_splits=5,shuffle=True,random_state=1)
optimal_params = RandomizedSearchCV(estimator=xgb.XGBRegressor(),
                                   param_distributions = param_grid, n_iter = 200,
                                   verbose = 1, scoring = 'neg_root_mean_squared_error',
                                   n_jobs=-1,random_state = 1,
                                   cv = cv)

optimal_params.fit(X,y)
print("Optimal parameter values =", optimal_params.best_params_)
print("Optimal cross validation R-squared = ",optimal_params.best_score_)
print("Time taken = ", round((time.time()-start_time)/60), " minutes")

```

Fitting 5 folds for each of 200 candidates, totalling 1000 fits

Optimal parameter values = {'subsample': 0.75, 'reg\_lambda': 1, 'n\_estimators': 1000, 'max\_d

Optimal cross validation R-squared = -5178.8689594137295

Time taken = 7 minutes

#RMSE based on the optimal parameter values

```
np.sqrt(mean_squared_error(optimal_params.best_estimator_.predict(Xtest),ytest))
```

5420.661056398766

## 11.2 CatBoost

CatBoost is a gradient boosting algorithm developed by Yandex (*Russian Google*) in 2017. Like LightGBM, CatBoost is also faster than XGBoost in training. However, unlike LightGBM, the authors have claimed that it outperforms both LightGBM and XGBoost in terms of prediction accuracy as well.

The key feature of CatBoost that address the issue with the gradient boosting procedure is the idea of ordered boosting. Classic boosting algorithms are prone to overfitting on small/noisy datasets due to a problem known as prediction shift. Recall, in gradient boosting, we fit trees on the gradient of the loss function (*refer the gradient boosting algorithm in section 10.10.2 of the book, [Elements of Statistical Learning](#)*):

$$r_m = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}.$$

When calculating the gradient estimate of an observation, these algorithms use the same observations that the model was built with, thus having no chances of experiencing unseen

data. CatBoost, on the other hand, uses the concept of ordered boosting, a permutation-driven approach to train model on a subset of data while calculating residuals on another subset, thus preventing “target leakage” and overfitting. The residuals of an observation are computed based on a model developed on the previous observations, where the observations are randomly shuffled at each iteration, i.e., for each tree.

Thus, the gradient of the loss function is based on test (*unseen*) data, instead of the data on which the model has been trained, which improves the generalizability of the model, and avoids overfitting on train data.

The authors have also shown that CatBoost performs better than XGBoost and LightGBM without tuning, i.e., with default hyperparameter settings.

Read the [CatBoost paper](#) for more details.

Here is a good [blog](#) listing the key features of CatBoost.

### 11.2.1 CatBoost for regression

We’ll use the function [CatBoostRegressor](#) for regression. For classification problems [CatBoostClassifier](#) can be used.

Let us check the performance of `CatBoostRegressor()` without tuning, i.e., with default hyperparameter settings.

```
model_cat = CatBoostRegressor(verbose=0).fit(X, y)

cv = KFold(n_splits=5, shuffle=True, random_state=1)
np.mean(-cross_val_score(CatBoostRegressor(verbose=0), X, y, cv = cv, n_jobs = -1,
                          scoring='neg_root_mean_squared_error'))
```

```
5035.972129299527
```

```
np.sqrt(mean_squared_error(model_cat.predict(Xtest), ytest))
```

```
5288.82153844634
```

Even with default hyperparameter settings, CatBoost has outperformed both XGBoost and LightGBM in terms of cross-validated RMSE, and RMSE on test data for our example of predicting car prices.

### 11.2.2 Target encoding with CatBoost

Target encoding for categorical variables can be used with CatBoost, that may benefit in terms of both speed and accuracy. However, the benefit is not guaranteed. Let us use target encoding for the categorical predictors `brand`, `model`, `transmission` and `fuelType`.

```
X = train[['mileage','mpg','year','engineSize', 'brand', 'model', 'transmission', 'fuelType']]
Xtest = test[['mileage','mpg','year','engineSize', 'brand', 'model', 'transmission', 'fuelType']]
y = train['price']
ytest = test['price']
```

The parameter `cat_features` will be used to specify the indices of the categorical predictors for target encoding.

```
model = CatBoostRegressor(verbose = False, cat_features = range(4, 8)).fit(X, y)
mean_squared_error(model.predict(Xtest), ytest, squared = False)
```

3263.1348853593345

Let us compare the results with one-hot encoding of the categorical predictors.

```
X = train[['mileage','mpg','year','engineSize', 'brand', 'model', 'transmission', 'fuelType']]
Xtest = test[['mileage','mpg','year','engineSize', 'brand', 'model', 'transmission', 'fuelType']]
X = pd.get_dummies(X)
Xtest = pd.get_dummies(Xtest)
```

In one-hot encoding, we need to make sure that both the datasets have the same predictors. Let us find the predictors in train data that are not in test data. Note that this is not necessary in target encoding.

```
np.setdiff1d(X.columns, Xtest.columns)
```

```
array(['model_ M6'], dtype=object)
```

```
X.drop(columns = 'model_ M6', inplace = True)
y = train['price']
ytest = test['price']
```

```
model = CatBoostRegressor(verbose = False).fit(X, y)
mean_squared_error(model.predict(Xtest), ytest, squared = False)
```

3219.857899121199

In this case, target encoding has a slightly higher RMSE as compared to one-hot encoding. However, it may do better than one-hot-encoding in a different problem.

Let us use both target encoding and one-hot encoding together to see if it helps do better than each of the them individually.

```
X = pd.concat([train[['brand', 'model', 'transmission', 'fuelType']], X], axis = 1)
Xtest = pd.concat([test[['brand', 'model', 'transmission', 'fuelType']], Xtest], axis = 1)
```

```
model = CatBoostRegressor(verbose = False, cat_features = range(4)).fit(X, y)
mean_squared_error(model.predict(Xtest), ytest, squared = False)
```

3172.449374536484

In this case, using target-encoding and one-hot-encoding together does better on test data. Using both the encodings together will help reduce bias while increasing variance. The benefit of using both the encodings together depends on the bias-variance tradeoff.

### 11.2.3 CatBoost vs XGBoost

Let us see the performance of XGBoost with default hyperparameter settings.

```
model_xgb = xgb.XGBRFRegressor().fit(X, y)
np.mean(-cross_val_score(xgb.XGBRFRegressor(), X, y, cv = cv, n_jobs = -1,
                        scoring='neg_root_mean_squared_error'))
```

6273.043859096154

```
np.sqrt(mean_squared_error(model_xgb.predict(Xtest), ytest))
```

6821.745153860935

XGBoost performance deteriorates showing that hyperparameter tuning is more important in XGBoost.

Let us see the performance of LightGBM with default hyperparameter settings.

```
model_lgbm = LGBMRegressor().fit(X, y)
np.mean(-cross_val_score(LGBMRegressor(), X, y, cv = cv, n_jobs = -1,
                          scoring='neg_root_mean_squared_error'))
```

5562.149251902867

```
np.sqrt(mean_squared_error(model_lgbm.predict(Xtest),ytest))
```

5494.0777923513515

LightGBM's default hyperparameter settings also seem to be more robust as compared to those of XGBoost.

### 11.2.4 Tuning CatBoostRegressor

The CatBoost hyperparameters can be tuned just like the XGBoost hyperparameters. However, there is some difference in the hyperparameters of both the packages. For example, `reg_alpha` (the *L1 penalization on weights of leaves*) and `colsample_bytree` (*subsample ratio of columns when constructing each tree*) hyperparameters are not there in CatBoost.

```
#K-fold cross validation to find optimal parameters for CatBoost regressor
start_time = time.time()
param_grid = {'max_depth': [4,6,8, 10],
              'num_leaves': [20, 31, 40, 60],
              'learning_rate': [0.01, 0.05, 0.1],
              'reg_lambda': [0, 10, 100],
              'n_estimators': [500, 1000, 1500],
              'subsample': [0.5, 0.75, 1.0],
              'colsample_bylevel': [0.25, 0.5, 0.75, 1.0]}

cv = KFold(n_splits=5,shuffle=True,random_state=1)
optimal_params = RandomizedSearchCV(estimator=CatBoostRegressor(random_state=1, verbose=False),
                                   grow_policy='Lossguide'),
                                   param_distributions = param_grid, n_iter = 200,
                                   verbose = 1,random_state = 1, scoring='neg_root_mean_squared_er
```

```

n_jobs=-1,
cv = cv)

optimal_params.fit(X,y)
print("Optimal parameter values =", optimal_params.best_params_)
print("Optimal cross validation RMSE = ",optimal_params.best_score_)
print("Time taken = ", round((time.time()-start_time)/60), " minutes")

```

Fitting 5 folds for each of 200 candidates, totalling 1000 fits

Optimal parameter values = {'subsample': 0.5, 'reg\_lambda': 0, 'num\_leaves': 40, 'n\_estimators': 2000}

Optimal cross validation RMSE = -4993.129407810791

Time taken = 23 minutes

```
#RMSE based on the optimal parameter values
```

```
np.sqrt(mean_squared_error(optimal_params.best_estimator_.predict(Xtest),ytest))
```

5249.434282204398

It takes 2 minutes to tune CatBoost, which is higher than LightGBM and lesser than XGBoost. CatBoost falls in between LightGBM and XGBoost in terms of speed. However, it is likely to be more accurate than XGBoost and LightGBM, and likely to require lesser tuning as compared to XGBoost.

```
model = CatBoostRegressor(grow_policy='Lossguide')
```

```
grid = {'num_leaves': Integer(4, 64),
        'learning_rate': Real(0.0001, 1.0),
        'reg_lambda': Real(0, 1e4),
        'n_estimators': Integer(2, 2000),
        'subsample': Real(0.1, 1.0),
        'colsample_bylevel': Real(0.1, 1.0)}
```

```
kfold = KFold(n_splits = 5, shuffle = True, random_state = 1)
```

```
gcv = BayesSearchCV(model, search_spaces = grid, cv = kfold, n_iter = 200, random_state = 1,
                    scoring = 'neg_root_mean_squared_error', n_jobs = -1)
```

```
paras = list(gcv.search_spaces.keys())
```

```
paras.sort()
```

```
def monitor(optim_result):
```

```
    cv_values = pd.Series(optim_result['func_vals']).cummin()
```

```
    display.clear_output(wait = True)
```

```

min_ind = pd.Series(optim_result['func_vals']).argmin()
print(paras, "=", optim_result['x_iters'][min_ind], pd.Series(optim_result['func_vals'])
sns.lineplot(cv_values)
plt.show()
gcv.fit(X, y, callback = monitor)

```

<IPython.core.display.Image object>

```

# Optimal values obtained
#['colsample_bylevel', 'learning_rate', 'n_estimators', 'num_leaves', 'reg_lambda', 'subsamp
#[0.3745508446405472, 0.1000958551500621, 2000, 11, 0.0, 0.3877212027881348] 5132.5378396768

```

### 11.2.5 Tuning Tips

Check the [documentation](#) for some tuning tips.

1. It is not recommended to use values greater than 64 for `num_leaves`, since it can significantly slow down the training process.
2. In most cases, the optimal depth ranges from 4 to 10. Values in the range from 6 to 10 are recommended. The maximum possible value of `max_depth` is 16.
3. Do not use one-hot encoding during preprocessing. This affects both the training speed and the resulting quality.
4. Symmetric trees have a very good prediction speed (roughly 10 times faster than non-symmetric trees) and give better quality in many cases.

## 12 Ensemble modeling

Ensembling models can help reduce error by leveraging the diversity and collective wisdom of multiple models. When ensembling, several individual models are trained independently and their predictions are combined to make the final prediction.

We have already seen examples of ensemble models in chapters 5 - 13. The ensembled models may reduce error by reducing the bias (*boosting*) and / or reducing the variance (*bagging* / *random forests* / *boosting*).

However, in this chapter we'll ensemble different types of models, instead of the same type of model. We may ensemble a linear regression model, a random forest, a gradient boosting model, and as many different types of models as we wish.

Below are a couple of reasons why ensembling models can be effective in reducing error:

1. **Bias reduction:** Different models may have different biases and the ensemble can help mitigate the individual biases, leading to a more generalized and accurate prediction. For example, consider that one model has a positive bias, and another model has a negative bias for the same instance. By averaging or combining the predictions of the two models, the biases may cancel out.
2. **Variance reduction:** As seen in the case of random forests and bagged trees, by averaging or combining the predictions of multiple models, the ensemble can reduce the overall variance and improve the accuracy of the final prediction. Note that for variance reduction, the models should have a low correlation (*recall the variance reduction formula of random forests*).

Mathematically also, we can show the effectiveness of an ensemble model. Let's consider the case of regression, and let the predictors be denoted as  $X$ , and the response as  $Y$ . Let  $f_1, \dots, f_m$  be individual models. The expected MSE of an ensemble can be written as:

$$E(MSE_{Ensemble}) = E\left[\left(\frac{1}{m} \sum_{i=1}^m f_i(X) - Y\right)^2\right] = \frac{1}{m^2} \sum_{i=1}^m E\left[(f_i(X) - Y)^2\right] + \frac{1}{m^2} \sum_{i \neq j} E\left[(f_i(X) - Y)(f_j(X) - Y)\right]$$

$$\Rightarrow E(MSE_{Ensemble}) = \frac{1}{m} \left( \frac{1}{m} \sum_{i=1}^m E\left[(f_i(X) - Y)^2\right] \right) + \frac{1}{m^2} \sum_{i \neq j} E\left[(f_i(X) - Y)(f_j(X) - Y)\right]$$



$$\Rightarrow E(MSE_{Ensemble}) = \frac{1}{m} \left( \frac{1}{m} \sum_{i=1}^m E(MSE_{f_i}) \right) + \frac{1}{m^2} \sum_{i \neq j} E \left[ (f_i(X) - Y)(f_j(X) - Y) \right]$$

If  $f_1, \dots, f_m$  are unbiased, then,

$$E(MSE_{Ensemble}) = \frac{1}{m} \left( \frac{1}{m} \sum_{i=1}^m E(MSE_{f_i}) \right) + \frac{1}{m^2} \sum_{i \neq j} Cov(f_i(X), f_j(X))$$

Assuming the **models are uncorrelated** (*i.e.*, they have a zero correlation), the second term (covariance of  $f_i(\cdot)$  and  $f_j(\cdot)$ ) reduces to zero, and the expected MSE of the ensemble reduces to:

$$E(MSE_{Ensemble}) = \frac{1}{m} \left( \frac{1}{m} \sum_{i=1}^m E(MSE_{f_i}) \right) \quad (12.1)$$

Thus, the expected MSE of an ensemble model with uncorrelated models is much smaller than the average MSE of all the models. Unless there is a model that is much better than the rest of the models, the MSE of the ensemble model is likely to be lower than the MSE of the individual models. However, there is no guarantee that the MSE of the ensemble model will be lower than the MSE of the individual models. Consider an extreme case where only one of the models have a zero MSE. The MSE of this model will be lower than the expected MSE of the ensemble model.

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import cross_val_score, train_test_split, GridSearchCV, ParameterGrid, StratifiedKFold, RandomizedSearchCV
from sklearn.metrics import mean_squared_error, r2_score, roc_curve, auc, precision_recall_curve
from sklearn.model_selection import KFold
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
from sklearn.ensemble import VotingRegressor, VotingClassifier, StackingRegressor, \
StackingClassifier, GradientBoostingRegressor, GradientBoostingClassifier, BaggingRegressor, \
BaggingClassifier, RandomForestRegressor, RandomForestClassifier, AdaBoostRegressor, AdaBoostClassifier
from sklearn.linear_model import LinearRegression, LogisticRegression, LassoCV, RidgeCV, ElasticNetCV
from sklearn.neighbors import KNeighborsRegressor
import itertools as it
```

```
import time as time
import xgboost as xgb
from catboost import CatBoostRegressor
from lightgbm import LGBMRegressor
from sklearn.preprocessing import PolynomialFeatures
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
```

```
#Using the same datasets as used for linear regression in STAT303-2,
#so that we can compare the non-linear models with linear regression
trainf = pd.read_csv('./Datasets/Car_features_train.csv')
trainp = pd.read_csv('./Datasets/Car_prices_train.csv')
testf = pd.read_csv('./Datasets/Car_features_test.csv')
testp = pd.read_csv('./Datasets/Car_prices_test.csv')
train = pd.merge(trainf, trainp)
test = pd.merge(testf, testp)
train.head()
```

	carID	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	18473	bmw	6 Series	2020	Semi-Auto	11	Diesel	145	53.3282	3.0	37980
1	15064	bmw	6 Series	2019	Semi-Auto	10813	Diesel	145	53.0430	3.0	33980
2	18268	bmw	6 Series	2020	Semi-Auto	6	Diesel	145	53.4379	3.0	36850
3	18480	bmw	6 Series	2017	Semi-Auto	18895	Diesel	145	51.5140	3.0	25998
4	18492	bmw	6 Series	2015	Automatic	62953	Diesel	160	51.4903	3.0	18990

```
X = train[['mileage', 'mpg', 'year', 'engineSize']]
Xtest = test[['mileage', 'mpg', 'year', 'engineSize']]
y = train['price']
ytest = test['price']
```

## 12.1 Ensembling regression models

### 12.1.1 Voting Regressor

Here, we will combine the predictions of different models. The function `VotingRegressor()` averages the predictions of all the models.

Below are the individual models tuned in the previous chapters.

```

#Tuned AdaBoost model from Section 7.2.4
model_ada = AdaBoostRegressor(estimator=DecisionTreeRegressor(max_depth=10),n_estimators=50,
                             learning_rate=1.0, random_state=1).fit(X, y)
print("RMSE for AdaBoost = ", np.sqrt(mean_squared_error(model_ada.predict(Xtest), ytest)))

#Tuned Random forest model from Section 6.1.2
model_rf = RandomForestRegressor(n_estimators=300, random_state=1,
                                n_jobs=-1, max_features=2).fit(X, y)
print("RMSE for Random forest = ", np.sqrt(mean_squared_error(model_rf.predict(Xtest), ytest)))

# Tuned XGBoost model from Section 9.2.6
model_xgb = xgb.XGBRegressor(random_state=1,max_depth=8,n_estimators=1000, subsample = 0.75,
                             learning_rate = 0.01,reg_lambda=1, gamma = 100).fit(X, y)
print("RMSE for XGBoost = ", np.sqrt(mean_squared_error(model_xgb.predict(Xtest), ytest)))

#Tuned gradient boosting model from Section 8.2.5
model_gb = GradientBoostingRegressor(max_depth=8,n_estimators=100,learning_rate=0.1,
                                     random_state=1,loss='huber').fit(X, y)
print("RMSE for Gradient Boosting = ", np.sqrt(mean_squared_error(model_gb.predict(Xtest), ytest)))

# Tuned Light GBM model from Section 13.1.1
model_lgbm = LGBMRegressor(subsample = 0.5, reg_lambda = 0, reg_alpha = 100, boosting_type =
                            num_leaves = 31, n_estimators = 500, learning_rate = 0.05, colsample_bytree = 1.0,
                            top_rate = 0.5).fit(X, y)
print("RMSE for LightGBM = ", np.sqrt(mean_squared_error(model_lgbm.predict(Xtest), ytest)))

# Tuned CatBoost model from Section 13.2.3
model_cat = CatBoostRegressor(subsample=0.5, num_leaves=40, n_estimators=500, max_depth=10,
                              verbose = False, learning_rate = 0.05, colsample_bylevel=0.75,
                              grow_policy='Lossguide', random_state = 1).fit(X, y)
print("RMSE for CatBoost = ", np.sqrt(mean_squared_error(model_cat.predict(Xtest), ytest)))

```

```

RMSE for AdaBoost = 5693.165811600585
RMSE for Random forest = 5642.45839697972
RMSE for XGBoost = 5497.553788113875
RMSE for Gradient Boosting = 5405.787029062213
RMSE for LightGBM = 5355.964600884197
RMSE for CatBoost = 5271.104736146779

```

Note that we **don't need to fit** the models **individually** before fitting them simultaneously in the voting ensemble. If we fit them individual, it will unnecessarily **waste time**.

Let us ensemble the models using the voting ensemble with equal weights.

```
#Voting ensemble: Averaging the predictions of all models

#Tuned AdaBoost model from Section 7.2.4
model_ada = AdaBoostRegressor(estimator=DecisionTreeRegressor(max_depth=10),
                              n_estimators=50, learning_rate=1.0, random_state=1)

#Tuned Random forest model from Section 6.1.2
model_rf = RandomForestRegressor(n_estimators=300, random_state=1,
                                n_jobs=-1, max_features=2)

# Tuned XGBoost model from Section 9.2.6
model_xgb = xgb.XGBRegressor(random_state=1, max_depth=8, n_estimators=1000, subsample = 0.75,
                              colsample_bytree = 1.0, learning_rate = 0.01, reg_lambda=1, gamma = 100)

#Tuned gradient boosting model from Section 8.2.5
model_gb = GradientBoostingRegressor(max_depth=8, n_estimators=100, learning_rate=0.1,
                                     random_state=1, loss='huber')

# Tuned CatBoost model from Section 13.2.3
model_cat = CatBoostRegressor(subsample=0.5, num_leaves=40, n_estimators=500, max_depth=10,
                              learning_rate = 0.05, colsample_bylevel=0.75, grow_policy='Loss',
                              random_state=1, verbose = False)

# Tuned Light GBM model from Section 13.1.1
model_lgbm = LGBMRegressor(subsample = 0.5, reg_lambda = 0, reg_alpha = 100, boosting_type =
                           num_leaves = 31, n_estimators = 500, learning_rate = 0.05,
                           colsample_bytree = 1.0, top_rate = 0.5)

start_time = time.time()
en = VotingRegressor(estimators = [('xgb', model_xgb), ('ada', model_ada), ('rf', model_rf),
                                  ('gb', model_gb), ('cat', model_cat), ('lgbm', model_lgbm)], n_jobs = -1)
en.fit(X, y)
print("Ensemble model RMSE = ", np.sqrt(mean_squared_error(en.predict(Xtest), ytest)))
print("Time taken = ", np.round((time.time() - start_time)/60, 2), "minutes")
```

Ensemble model RMSE = 5259.899392611916

Time taken = 0.21 minutes

As expected, RMSE of the ensembled model is less than that of each of the individual models.

Note that the RMSE can be **further improved** by **removing** the **weaker models** from the ensemble. Let us remove the three weakest models - XGBoost, Random forest, and Adaboost.

```
#Voting ensemble: Averaging the predictions of all models

start_time = time.time()
en = VotingRegressor(estimators = [('gb', model_gb), ('cat', model_cat), ('lgbm', model_lgbm)])
en.fit(X,y)
print("Ensemble model RMSE = ", np.sqrt(mean_squared_error(en.predict(Xtest), ytest)))
print("Time taken = ", np.round((time.time() - start_time)/60, 2), "minutes")
```

```
Ensemble model RMSE = 5191.814866810768
Time taken = 0.18 minutes
```

### 12.1.2 Stacking Regressor

Stacking is a more sophisticated method of ensembling models. The method is as follows:

1. The training data is split into  $K$  folds. Each of the  $K$  folds serves as a test data in one of the  $K$  iterations, and the rest of the folds serve as train data.
2. Each model is used to make predictions on each of the  $K$  folds, after being trained on the remaining  $K-1$  folds. In this manner, each model predicts the response on each train data point - when that train data point was not used to train the model.
3. Predictions at each training data points are generated by each model in step 2 (the above step). These predictions are now used as predictors to train a meta-model (referred by the argument `final_estimator`), with the original response as the response. The meta-model (or `final_estimator`) learns to combine predictions of different models to make a better prediction.

#### 12.1.2.1 Metamodel: Linear regression

```
#Stacking using LinearRegression as the metamodel
en = StackingRegressor(estimators = [('xgb', model_xgb), ('ada', model_ada), ('rf', model_rf),
                                   ('gb', model_gb), ('cat', model_cat), ('lgbm', model_lgbm)],
                      final_estimator=LinearRegression(),
                      cv = KFold(n_splits = 5, shuffle = True, random_state=1))
start_time = time.time()
en.fit(X,y)
```

```
print("Linear regression metamodel RMSE = ", np.sqrt(mean_squared_error(en.predict(Xtest), ytest)))
print("Time taken = ", np.round((time.time() - start_time)/60,2), "minutes")
```

```
Linear regression metamodel RMSE = 5220.456280327686
Time taken = 2.03 minutes
```

```
#Co-efficients of the meta-model
en.final_estimator_.coef_
```

```
array([ 0.05502964,  0.14566665,  0.01093624,  0.30478283,  0.57403909,
        -0.07057344])
```

```
sum(en.final_estimator_.coef_)
```

```
1.0198810182715363
```

Note the above coefficients of the meta-model. The model gives the **highest weight** to the **gradient boosting** model (*with **huber** loss*), and the **catboost** model, and the **lowest weight** to the relatively weak **random forest** model.

Also, note that the **coefficients need not sum to one**.

Let us try improving the RMSE further by removing the weaker models from the ensemble. Let us remove the three weakest models based on the size of their coefficients in the linear regression metamodel.

```
#Stacking using LinearRegression as the metamodel
en = StackingRegressor(estimators = [('gb', model_gb), ('cat', model_cat), ('ada', model_ada)],
                       final_estimator=LinearRegression(),
                       cv = KFold(n_splits = 5, shuffle = True, random_state=1))
start_time = time.time()
en.fit(X,y)
print("Linear regression metamodel RMSE = ", np.sqrt(mean_squared_error(en.predict(Xtest), ytest)))
print("Time taken = ", np.round((time.time() - start_time)/60,2), "minutes")
```

```
Linear regression metamodel RMSE = 5205.225710180056
Time taken = 1.36 minutes
```

The metamodel accuracy **improves further**, when **strong models** are ensembled.

```
#Co-efficients of the meta-model
en.final_estimator_.coef_
```

```
array([0.31824119, 0.54231032, 0.15998634])
```

```
sum(en.final_estimator_.coef_)
```

```
1.020537847948332
```

### 12.1.2.2 Metamodel: Lasso

```
#Stacking using Lasso as the metamodel
en = StackingRegressor(estimators = [('xgb', model_xgb), ('ada', model_ada), ('rf', model_rf),
                                   ('gb', model_gb), ('cat', model_cat), ('lgbm', model_lgbm) ],
                      final_estimator = LassoCV(),
                      cv = KFold(n_splits = 5, shuffle = True, random_state=1))
start_time = time.time()
en.fit(X,y)
print("Lasso metamodel RMSE = ", np.sqrt(mean_squared_error(en.predict(Xtest), ytest)))
print("Time taken = ", np.round((time.time() - start_time)/60,2), "minutes")
```

```
Lasso metamodel RMSE = 5206.021083501416
Time taken = 2.05 minutes
```

```
#Coefficients of the lasso metamodel
en.final_estimator_.coef_
```

```
array([ 0.03524446,  0.15077605,  0.          ,  0.30392268,  0.52946243,
        -0.          ])
```

Note that lasso **reduces the weight** of the **weak random forest** model, and **light gbm** model to **0**. Even though light GBM is a strong model, it may be **correlated or collinear** with XGBoost, or other models, and hence is not needed.

Note that as lasso performs **model selection** on its own, removing models with zero coefficients or weights does not make a difference, as shown below.

```
#Stacking using Lasso as the metamodel
en = StackingRegressor(estimators = [('xgb', model_xgb), ('ada', model_ada),
                                   ('gb', model_gb), ('cat', model_cat) ],
                      final_estimator = LassoCV(),
                      cv = KFold(n_splits = 5, shuffle = True, random_state=1))
start_time = time.time()
en.fit(X,y)
print("Lasso metamodel RMSE = ", np.sqrt(mean_squared_error(en.predict(Xtest),ytest)))
print("Time taken = ", np.round((time.time() - start_time)/60,2), "minutes")
```

```
Lasso metamodel RMSE = 5205.93233977352
Time taken = 1.79 minutes
```

```
#Coefficients of the lasso metamodel
en.final_estimator_.coef_
```

```
array([0.03415944, 0.15053122, 0.30464838, 0.53006297])
```

### 12.1.2.3 Metamodel: Random forest

A highly flexible model such as a random forest may not be a good choice for ensembling correlated models. However, let us tune the random forest meta model, and check its accuracy.

```
# Tuning hyperparameter of the random forest meta-model
start_time = time.time()
oob_score_i = []
for i in range(1, 7):
    en = StackingRegressor(estimators = [('xgb', model_xgb), ('ada', model_ada), ('rf', model_rf),
                                       ('gb', model_gb), ('cat', model_cat), ('lgbm', model_lgbm)],
                          final_estimator = RandomForestRegressor(max_features = i, oob_score = True),
                          cv = KFold(n_splits = 5, shuffle = True, random_state=1)).fit(X,y)
    oob_score_i.append(en.final_estimator_.oob_score_)
print("Time taken = ", np.round((time.time() - start_time)/60,2), "minutes")
```

```
Time taken = 12.08 minutes
```

```
print("Optimal value of max_features =", np.array(oob_score_i).argmax() + 1)
```

```
Optimal value of max_features = 1
```



```
# Training the tuned random forest metamodel
start_time = time.time()
en = StackingRegressor(estimators = [('xgb', model_xgb), ('ada', model_ada),
                                   ('rf', model_rf), ('gb', model_gb), ('cat', model_cat),
                                   ('lgbm', model_lgbm)],
                      final_estimator = RandomForestRegressor(max_features = 1,
                                                              n_estimators=500), cv = KFold(n_splits = 5, shuffle = True,
                                                              random_state=1)).fit(X,y)
print("Random Forest metamodel RMSE = ", np.sqrt(mean_squared_error(en.predict(Xtest),ytest)))
print("Time taken = ", np.round((time.time() - start_time)/60,2), "minutes")
```

```
Random Forest metamodel RMSE = 5441.9155087961
Time taken = 1.71 minutes
```

Note that highly flexible models may not be needed when the predictors are highly correlated with the response. However, in some cases, they may be useful, as in the classification example in the next section.

#### 12.1.2.4 Metamodel: CatBoost

```
#Stacking using MARS as the meta-model
en = StackingRegressor(estimators = [('xgb', model_xgb), ('ada', model_ada), ('rf', model_rf),
                                   ('gb', model_gb), ('cat', model_cat), ('lgbm', model_lgbm)],
                      final_estimator = CatBoostRegressor(verbose = False),
                      cv = KFold(n_splits = 5, shuffle = True, random_state=1))
start_time = time.time()
en.fit(X,y)
print("Random Forest metamodel RMSE = ", np.sqrt(mean_squared_error(en.predict(Xtest),ytest)))
print("Time taken = ", np.round((time.time() - start_time)/60,2), "minutes")
```

```
Random Forest metamodel RMSE = 5828.803609683251
Time taken = 1.66 minutes
```

## 12.2 Ensembling classification models

We'll ensemble models for predicting accuracy of identifying people having a heart disease.

```

data = pd.read_csv('./Datasets/Heart.csv')
data.dropna(inplace = True)
#Response variable
y = pd.get_dummies(data['AHD'])['Yes']

#Creating a dataframe for predictors with dummy variables replacing the categorical variables
X = data.drop(columns = ['AHD','ChestPain','Thal'])
X = pd.concat([X,pd.get_dummies(data['ChestPain']),pd.get_dummies(data['Thal'])],axis=1)

#Creating train and test datasets
Xtrain,Xtest,ytrain,ytest = train_test_split(X,y,train_size = 0.5,random_state=1)

```

Let us tune the individual models first.

## AdaBoost

```

# Tuning Adaboost for maximizing accuracy
model = AdaBoostClassifier(random_state=1)
grid = dict()
grid['n_estimators'] = [10, 50, 100,200,500]
grid['learning_rate'] = [0.0001, 0.001, 0.01,0.1, 1.0]
grid['base_estimator'] = [DecisionTreeClassifier(max_depth=1), DecisionTreeClassifier(max_depth=2),
                           DecisionTreeClassifier(max_depth=3),DecisionTreeClassifier(max_depth=4)]
# define the evaluation procedure
cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=1)
# define the grid search procedure
grid_search = GridSearchCV(estimator=model, param_grid=grid, n_jobs=-1, cv=cv, scoring='accuracy')
# execute the grid search
grid_result = grid_search.fit(Xtrain, ytrain)
# summarize the best score and configuration
print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))

```

Best: 0.871494 using {'base\_estimator': DecisionTreeClassifier(max\_depth=1), 'learning\_rate': 0.0001}

## Gradient Boosting

```

# Tuning gradient boosting for maximizing accuracy
model = GradientBoostingClassifier(random_state=1)
grid = dict()
grid['n_estimators'] = [10, 50, 100,200,500]
grid['learning_rate'] = [0.0001, 0.001, 0.01,0.1, 1.0]
grid['max_depth'] = [1,2,3,4,5]
grid['subsample'] = [0.5,1.0]
# define the evaluation procedure
cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=1)
# define the grid search procedure
grid_search = GridSearchCV(estimator=model, param_grid=grid, n_jobs=-1, cv=cv, scoring='accuracy')
# execute the grid search
grid_result = grid_search.fit(Xtrain, ytrain)
# summarize the best score and configuration
print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))

```

Best: 0.871954 using {'learning\_rate': 1.0, 'max\_depth': 4, 'n\_estimators': 100, 'subsample': 0.5}

## XGBoost

```

# Tuning XGBoost for maximizing accuracy
start_time = time.time()
param_grid = {'n_estimators': [25, 100,250,500],
              'max_depth': [4, 6 ,8],
              'learning_rate': [0.01,0.1,0.2],
              'gamma': [0, 1, 10, 100],
              'reg_lambda': [0, 10, 100],
              'subsample': [0.5, 0.75, 1.0]
              'scale_pos_weight': [1.25,1.5,1.75]}#Control the balance of positive and negative samples

cv = StratifiedKFold(n_splits=5,shuffle=True,random_state=1)
optimal_params = GridSearchCV(estimator=xgb.XGBClassifier(random_state=1),
                              param_grid = param_grid,
                              scoring = 'accuracy',
                              verbose = 1,
                              n_jobs=-1,
                              cv = cv)
optimal_params.fit(Xtrain,ytrain)

```

```
print(optimal_params.best_params_,optimal_params.best_score_)
print("Time taken = ", (time.time()-start_time)/60, " minutes")
```

Fitting 5 folds for each of 972 candidates, totalling 4860 fits

```
{'gamma': 0, 'learning_rate': 0.2, 'max_depth': 4, 'n_estimators': 25, 'reg_lambda': 0, 'sca
```

Time taken = 0.9524135629336039 minutes

```
#Tuned Adaboost model
```

```
model_ada = AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_depth=1), n_estimators=500,
                               random_state=1, learning_rate=0.01).fit(Xtrain, ytrain)
```

```
test_accuracy_ada = model_ada.score(Xtest,ytest) #Returns the classification accuracy of the model
```

```
#Tuned Random forest model from Section 6.3
```

```
model_rf = RandomForestClassifier(n_estimators=500, random_state=1,max_features=3,
                                 n_jobs=-1,oob_score=False).fit(Xtrain, ytrain)
```

```
test_accuracy_rf = model_rf.score(Xtest,ytest) #Returns the classification accuracy of the model
```

```
#Tuned gradient boosting model
```

```
model_gb = GradientBoostingClassifier(n_estimators=100, random_state=1,max_depth=4,learning_rate=0.1,
                                     subsample = 1.0).fit(Xtrain, ytrain)
```

```
test_accuracy_gb = model_gb.score(Xtest,ytest) #Returns the classification accuracy of the model
```

```
#Tuned XGBoost model
```

```
model_xgb = xgb.XGBClassifier(random_state=1,gamma=0,learning_rate = 0.2,max_depth=4,
                              n_estimators = 25,reg_lambda = 0,scale_pos_weight=1.25).fit(Xtrain, ytrain)
```

```
test_accuracy_xgb = model_xgb.score(Xtest,ytest) #Returns the classification accuracy of the model
```

```
print("Adaboost accuracy = ",test_accuracy_ada)
```

```
print("Random forest accuracy = ",test_accuracy_rf)
```

```
print("Gradient boost accuracy = ",test_accuracy_gb)
```

```
print("XGBoost model accuracy = ",test_accuracy_xgb)
```

Adaboost accuracy = 0.7986577181208053

Random forest accuracy = 0.8120805369127517

Gradient boost accuracy = 0.7986577181208053

XGBoost model accuracy = 0.7785234899328859

### 12.2.1 Voting classifier - hard voting

In this type of ensembling, the predicted class is the one predicted by the majority of the classifiers.

```
ensemble_model = VotingClassifier(estimators=[('ada',model_ada),('rf',model_rf),('gb',model_gb)])
ensemble_model.fit(Xtrain,ytrain)
ensemble_model.score(Xtest, ytest)
```

0.825503355704698

Note that the prediction accuracy of the ensemble is higher than the prediction accuracy of each of the individual models on unseen data.

### 12.2.2 Voting classifier - soft voting

In this type of ensembling, the predicted class is the one based on the average predicted probabilities of all the classifiers. The threshold probability is 0.5.

```
ensemble_model = VotingClassifier(estimators=[('ada',model_ada),('rf',model_rf),('gb',model_gb)],
                                voting='soft')
ensemble_model.fit(Xtrain,ytrain)
ensemble_model.score(Xtest, ytest)
```

0.7919463087248322

Note that soft voting will be good only for well calibrated classifiers, i.e., all the classifiers must have probabilities at the same scale.

### 12.2.3 Stacking classifier

Conceptually, the idea is similar to that of Stacking regressor.

```
#Using Logistic regression as the meta model (final_estimator)
ensemble_model = StackingClassifier(estimators=[('ada',model_ada),('rf',model_rf),('gb',model_gb)],
                                   final_estimator=LogisticRegression(random_state=1,max_iter=1000),
                                   cv = StratifiedKFold(n_splits=5,shuffle=True,random_state=1))
ensemble_model.fit(Xtrain,ytrain)
ensemble_model.score(Xtest, ytest)
```

0.7986577181208053

```
#Coefficients of the logistic regression metamodel
ensemble_model.final_estimator_.coef_
```

```
array([[0.81748051, 1.28663164, 1.64593342, 1.50947087]])
```

```
#Using random forests as the meta model (final_estimator). Note that random forest will require tuning
ensemble_model = StackingClassifier(estimators=[('ada',model_ada),('rf',model_rf),('gb',model_gb)],
                                   final_estimator=RandomForestClassifier(n_estimators=500, n_jobs=-1,
                                                                           random_state=1,oob_score=True),
                                   cv = StratifiedKFold(n_splits=5,shuffle=True,random_state=1))
ensemble_model.fit(Xtrain,ytrain)
ensemble_model.score(Xtest, ytest)
```

```
0.8322147651006712
```

Note that a complex `final_estimator` such as random forest will require tuning. In the above case, the `max_features` argument of random forests has been tuned to obtain the maximum OOB score. The tuning is shown below.

```
#Tuning the random forest parameters
start_time = time.time()
oob_score = {}

i=0
for pr in range(1,5):
    model = StackingClassifier(estimators=[('ada',model_ada),('rf',model_rf),('gb',model_gb)],
                              final_estimator=RandomForestClassifier(n_estimators=500, n_jobs=-1,
                                                                      random_state=1,oob_score=True),
                              cv = StratifiedKFold(n_splits=5,shuffle=True,random_state=1))
    oob_score[pr] = model.final_estimator_.oob_score_

end_time = time.time()
print("time taken = ", (end_time-start_time)/60, " minutes")
print("max accuracy = ", np.max(list(oob_score.values())))
print("Best value of max_features= ", np.argmax(list(oob_score.values()))+1)
```

```
time taken = 0.33713538646698 minutes
max accuracy = 0.8445945945945946
Best value of max_features= 1
```

```
#The final predictor (metamodel) - random forest obtains the maximum oob_score for max_features
oob_score
```

```
{1: 0.8445945945945946,
 2: 0.831081081081081,
 3: 0.8378378378378378,
 4: 0.831081081081081}
```

## 12.2.4 Tuning all models simultaneously

Individual model hyperparameters can be tuned simultaneously while ensembling them with a `VotingClassifier()`. However, this approach can be too expensive for even moderately-sized datasets.

```
# Create the param grid with the names of the models as prefixes

model_ada = AdaBoostClassifier(base_estimator = DecisionTreeClassifier())
model_rf = RandomForestClassifier()
model_gb = GradientBoostingClassifier()
model_xgb = xgb.XGBClassifier()

ensemble_model = VotingClassifier(estimators=[('ada',model_ada),('rf',model_rf),('gb',model_gb)])

hp_grid = dict()

# XGBoost
hp_grid['xgb__n_estimators'] = [25, 100,250,50]
hp_grid['xgb__max_depth'] = [4, 6 ,8]
hp_grid['xgb__learning_rate'] = [0.01, 0.1, 1.0]
hp_grid['xgb__gamma'] = [0, 1, 10, 100]
hp_grid['xgb__reg_lambda'] = [0, 1, 10, 100]
hp_grid['xgb__subsample'] = [0, 1, 10, 100]
hp_grid['xgb__scale_pos_weight'] = [1.0, 1.25, 1.5]
hp_grid['xgb__colsample_bytree'] = [0.5, 0.75, 1.0]

# AdaBoost
hp_grid['ada__n_estimators'] = [10, 50, 100,200,500]
hp_grid['ada__base_estimator__max_depth'] = [1, 3, 5]
hp_grid['ada__learning_rate'] = [0.01, 0.1, 0.2]

# Random Forest
```

```

hp_grid['rf__n_estimators'] = [100]
hp_grid['rf__max_features'] = [3, 6, 9, 12, 15]

# GradBoost
hp_grid['gb__n_estimators'] = [10, 50, 100, 200, 500]
hp_grid['gb__max_depth'] = [1, 3, 5]
hp_grid['gb__learning_rate'] = [0.01, 0.1, 0.2, 1.0]
hp_grid['gb__subsample'] = [0.5, 0.75, 1.0]

start_time = time.time()
grid = RandomizedSearchCV(ensemble_model, hp_grid, cv=5, scoring='accuracy', verbose = True,
                          n_iter = 100, n_jobs=-1).fit(Xtrain, ytrain)
print("Time taken = ", round((time.time()-start_time)/60), " minutes")

grid.best_estimator_.score(Xtest, ytest)

```

0.8120805369127517

## 12.3 Ensembling models based on different sets of predictors

Generally, tree-based models such as CatBoost, and XGBoost are the most accurate, while other models, such as bagging, random forests, KNN, and linear models, may not be as accurate. Thus, sometimes, the weaker models, despite bringing-in diversity in the model ensemble may deteriorate the ensemble accuracy due to their poor individual performance (*check slides for technical details*). Thus, sometimes, another approach to bring-in model diversity is to develop strong models based on different sets of predictors, and ensemble them.

Different feature selection methods (*such as Lasso, feature importance returned by tree-based methods, stepwise k-fold feature selection, etc.*), may be used to obtain different sets of important features, strong models can be tuned on these sets, and then ensembled. Even though the models may be of the same type, the different sets of predictors will help bring-in the element of diversity in the ensemble.

```

trainf = pd.read_csv('./Datasets/Car_features_train.csv')
trainp = pd.read_csv('./Datasets/Car_prices_train.csv')
testf = pd.read_csv('./Datasets/Car_features_test.csv')
testp = pd.read_csv('./Datasets/Car_prices_test.csv')
train = pd.merge(trainf, trainp)
test = pd.merge(testf, testp)
train.head()

```



	carID	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	18473	bmw	6 Series	2020	Semi-Auto	11	Diesel	145	53.3282	3.0	37980
1	15064	bmw	6 Series	2019	Semi-Auto	10813	Diesel	145	53.0430	3.0	33980
2	18268	bmw	6 Series	2020	Semi-Auto	6	Diesel	145	53.4379	3.0	36850
3	18480	bmw	6 Series	2017	Semi-Auto	18895	Diesel	145	51.5140	3.0	25998
4	18492	bmw	6 Series	2015	Automatic	62953	Diesel	160	51.4903	3.0	18990

```
X = train[['mileage','mpg','year','engineSize']]
Xtest = test[['mileage','mpg','year','engineSize']]
y = train['price']
ytest = test['price']
```

We will create polynomial interactions to develop two sets of predictors - first order predictors, and second order predictors.

```
poly_set = PolynomialFeatures(2, include_bias = False)
X_poly = poly_set.fit_transform(X)
X_poly = pd.DataFrame(X_poly, columns=poly_set.get_feature_names_out())
X_poly.columns = X_poly.columns.str.replace("^", "_", regex=True)
Xtest_poly = poly_set.fit_transform(Xtest)
Xtest_poly = pd.DataFrame(Xtest_poly, columns=poly_set.get_feature_names_out())
Xtest_poly.columns = Xtest_poly.columns.str.replace("^", "_", regex=True)
```

Let us use 2 different sets of predictors to introduce diversity in the ensemble.

```
col_set1 = ['mileage','mpg', 'year','engineSize']
col_set2 = X_poly.columns
```

Let us use two types of strong tree-based models.

```
cat = CatBoostRegressor(verbose=False)
gb = GradientBoostingRegressor(loss = 'huber')
```

We will use the `Pipeline()` function along with `ColumnTransformer()` to map a predictor set to each model.

```
cat_pipe1 = Pipeline([
    ('column_transformer', ColumnTransformer([('cat1_transform', 'passthrough', col_set1)],
    ('cat1', cat)
```

```

])

cat_pipe2 = Pipeline([
    ('column_transformer', ColumnTransformer([('cat2_transform', 'passthrough', col_set2)], r
    ('cat2', cat)
])

gb_pipe1 = Pipeline([
    ('column_transformer', ColumnTransformer([('gb1_transform', 'passthrough', col_set1)], r
    ('gb1', gb)
])

gb_pipe2 = Pipeline([
    ('column_transformer', ColumnTransformer([('gb2_transform', 'passthrough', col_set2)], r
    ('gb2', gb)
])

```

We will use Linear regression to ensemble the models.

```
en_new.final_estimator_.coef_
```

```
array([ 0.30127482,  0.79242981, -0.07168258, -0.01781781])
```

```

en_new = StackingRegressor(estimators = [('cat1', cat_pipe1), ('cat2', cat_pipe2),
                                       ('gb1', gb_pipe1), ('gb2', gb_pipe2)],
                           final_estimator=LinearRegression(),
                           cv = KFold(n_splits = 15, shuffle = True, random_state=1))

```

```
en_new.fit(X_poly, y)
```

```

StackingRegressor(cv=KFold(n_splits=15, random_state=1, shuffle=True),
                  estimators=[('cat1',
                               Pipeline(steps=[('column_transformer',
                                                  ColumnTransformer(transformers=[('cat1_transf
                                                  'passthrough',
                                                  ['mileage',
                                                  'mpg',
                                                  'year',
                                                  'engineSize
                               ('cat1',

```

```

<catboost.core.CatBoostRegressor object at 0x...
('cat2',
 Pipeline(steps=[('column_transformer',...
 Pipeline(steps=[('column_transformer',
                  ColumnTransformer(transformers=[('gb2_transf
                  'passthrough',
                  Index(['mileage year', 'mileage engineSize', 'mpg_2', 'mpg year',
                  'mpg engineSize', 'year_2', 'year engineSize', 'engineSize_2'],
dtype='object'))])),
                  ('gb2',
                   GradientBoostingRegressor(loss='huber'))]))]
final_estimator=LinearRegression())

```

```
mean_squared_error(en_new.predict(Xtest_poly), ytest, squared = False)
```

```
5185.376722607323
```

Note that the above model does better on test data than all the models developed so far. Using different sets of predictors introduces diversity in the ensemble, as an alternative to including “weaker” models in the ensemble to add diversity.

Check the idea being used in the Spring 2023 prediction problem in the appendix.

# A Assignment 1

## Instructions

1. You may talk to a friend, discuss the questions and potential directions for solving them. However, you need to write your own solutions and code separately, and not as a group activity.
2. Write your code in the **Code cells** and your answers in the **Markdown cells** of the Jupyter notebook. Ensure that the solution is written neatly enough to for the graders to understand and follow.
3. Use [Quarto](#) to render the **.ipynb** file as HTML. You will need to open the command prompt, navigate to the directory containing the file, and use the command: `quarto render filename.ipynb --to html`. Submit the HTML file.
4. The assignment is worth 100 points, and is due on **Thursday, 11th April 2025 at 11:59 pm**.
5. **Five points are properly formatting the assignment.** The breakdown is as follows:
  - Must be an HTML file rendered using Quarto (**2 points**). *If you have a Quarto issue, you must mention the issue & quote the error you get when rendering using Quarto in the comments section of Canvas, and submit the ipynb file.*
  - There aren't excessively long outputs of extraneous information (e.g. no printouts of entire data frames without good reason, there aren't long printouts of which iteration a loop is on, there aren't long sections of commented-out code, etc.) (**1 point**)
  - Final answers to each question are written in the Markdown cells. (**1 point**)
  - There is no piece of unnecessary / redundant code, and no unnecessary / redundant text. (**1 point**)

## A.1 1) Bias-Variance Trade-off for Regression (50 points)

The main goal of this question is to understand and visualize the bias-variance trade-off in a regression model by performing repetitive simulations.

The conceptual clarity about bias and variance will help with the main logic behind creating many models that will come up later in the course.

### A.1.1 a) Define the True Relationship (Signal)

First, you need to implement the underlying true relationship (Signal) you want to sample data from. Assume that the function is the [Bukin function](#). Implement it as a user-defined function and run it with the test cases below to make sure it is implemented correctly. **(5 points)**

**Note:** It would be more useful to have only one input to the function. You can treat the input as an array of two elements.

```
print(Bukin(np.array([1,2]))) # The output should be 141.177
print(Bukin(np.array([6,-4]))) # The output should be 208.966
print(Bukin(np.array([0,1]))) # The output should be 100.1
```

### A.1.2 b) Generate Test Set (No Noise)

Generate a **noiseless** test set with **100 observations** sampled from the true underlying function. This test set will be used to evaluate **bias and variance**, so make sure it follows the correct data generation process.

**(5 points)**

**Instructions:**

- **Do not use loops** for this question.
- `.apply` will be especially helpful (and often simpler).

**Data generation assumptions:**

- Use `np.random.seed(100)` for reproducibility.
- The first predictor,  $x_1$ , should be drawn from a **uniform distribution** over the interval  $[-15, -5]$ , i.e.,  $x_1 \sim U[-15, -5]$ .
- The second predictor,  $x_2$ , should be drawn from a **uniform distribution** over the interval  $[-3, 3]$ , i.e.,  $x_2 \sim U[-3, 3]$ .
- Compute the true function values using the underlying model as your response  $y$

### A.1.3 c) Initialize Results DataFrame

Create an empty DataFrame with the following columns:

- **degree**: the degree of the polynomial model
- **bias\_sq**: estimated squared bias (averaged over test points)
- **var**: estimated variance of predictions
- **bias\_var\_sum**: sum of bias squared and variance
- **empirical\_mse**: mean squared error calculated using sklearn's `mean_squared_error()` on model predictions vs. true function values

This DataFrame will be used to store the results of your bias–variance tradeoff analysis and for generating comparison plots.

(3 points)

### A.1.4 d) Generate Training Sets (With Noise)

To estimate the **bias**, **variance**, and **total error (MSE)** of a Linear Regression model trained on noisy data from the underlying Bukin function, follow the steps below.

#### Step 1: Generate 100 Training Sets

- Create **100 independent training datasets**, each with **100 observations** (same size as the test set).
- For each training dataset:
  - Use `np.random.seed(i)` to ensure reproducibility, where `i` is the dataset index (0 to 99).
  - Sample predictors from the **same distributions** used to generate the test set.
  - Add **Gaussian noise** with mean 0 and standard deviation 10:  
 $\varepsilon \sim \mathcal{N}(0, 10)$

#### Step 2: Train Polynomial Models (Degrees 1 to 7)

- For each training dataset, train polynomial models with degrees **1 through 7**.
- Use polynomial feature transformations that include both:
  - **Higher-order terms** (e.g.,  $x_1^2$ ,  $x_2^3$ )
  - **Interaction terms** (e.g.,  $x_1 \cdot x_2$ )
- Make predictions on the **fixed, noiseless test set** for each trained model.

### Step 3: Estimate Bias<sup>2</sup>, Variance, and MSE

- For each **degree**, and each **test point**, collect the 100 predicted values from the models trained on the different training sets.
- Using these predictions, compute:
  - **Bias squared**: squared difference between the mean prediction and the true value.
  - **Variance**: variance of the predictions.
  - **Theoretical MSE**: sum of bias squared and variance.
  - **Empirical MSE**: compute using `sklearn.metrics.mean_squared_error` between each model's prediction and the true values, then average over the 100 training runs.
- Store all four quantities for each degree in your results DataFrame:
  - `degree`
  - `bias_sq`
  - `var`
  - `bias_var_sum` (bias squared + variance)
  - `empirical_mse`

(25 points)

#### Reminder: Comparing Theoretical vs. Empirical MSE

When evaluating model performance on the **noiseless test set**:

- The **irreducible error** (i.e., noise in training data) does **not** affect the test targets.
- Therefore, the test error (MSE) can be decomposed as:
$$MSE = Bias^2 + Variance$$
- The **empirical MSE** (from sklearn) should closely match the **sum of bias<sup>2</sup> and variance**, since the test data contains no noise.

### A.1.5 e) Visualize Bias–Variance Decomposition

Using the results stored in your DataFrame, create a plot with **four lines**, each plotted against the polynomial **degree**:

1. **Bias squared**
2. **Variance**
3. **Bias squared + Variance** (i.e., the theoretical decomposition of MSE)
4. **Empirical MSE** calculated using `sklearn.metrics.mean_squared_error()` (computed from the predicted values vs. true function values on the noiseless test set)

**Plot requirements:** - Use a single line plot with the polynomial degree on the x-axis and error values on the y-axis. - Include a **legend** to clearly label each line. - Use different line styles or markers for easy visual comparison.

**Goal:** - Compare the **empirical MSE** to the **sum of bias squared and variance**. - If everything is implemented correctly, the two lines should be very close (or even identical, up to numerical precision).

### A.1.6 f) Identify the Optimal Model

- What is the **optimal polynomial degree** based on the **lowest empirical MSE** (calculated using sklearn)?  
(2 points)
- Report the corresponding values of:
  - **Bias squared**
  - **Variance**
  - **Bias squared + Variance**
  - **Empirical MSE**  
for that degree.  
(3 points)

## A.2 2) Building a Low-Bias, Low-Variance Model via Regularization (50 points)

The main goal of this question is to further reduce the **total prediction error** by applying **regularization**.

Specifically, you'll use **Ridge regression** to build a **low-bias, low-variance** model for data generated from the underlying Bukin function with noise.

### A.2.1 a) Why Regularization?

Explain why the model with the optimal polynomial degree (as identified in Question 1) is **not guaranteed** to be the true low-bias, low-variance model.

Why might **regularization** still be necessary to improve generalization performance, even after selecting the degree that minimizes MSE?



(5 points)

### A.2.2 b) Which Degrees to Exclude?

Based on your plot and results from **1e** and **1f**, identify which polynomial degrees should be **excluded** from regularization experiments because they are already too simple (high bias) or too complex (high variance).

Explain which degrees you will exclude and **why**, using your understanding of how **regularization affects bias and variance**.

(10 points)

### A.2.3 c) Apply Ridge Regularization

Repeat the steps from **1c** and **1d**, but this time use **Ridge regression** instead of ordinary least squares.

- Use only the degrees **not excluded** in 2b (and also exclude degree 7 to avoid extreme overfitting).
- Use **5-fold cross-validation** to tune the Ridge regularization strength.
- Use `neg_root_mean_squared_error` as the scoring metric for cross-validation.
- Tune over a range of regularization strengths (e.g., from 1 to 100).
- For each retained degree, compute:
  - **Bias squared**
  - **Variance**
  - **Bias squared + Variance**
  - **Empirical MSE** (from `sklearn.metrics.mean_squared_error`)

Store your results in a new DataFrame with the same structure as in Question 1.

(10 points)

### A.2.4 d) Visualize Regularized Results

Repeat the visualization from **1e**, but using the results from **2c** (Ridge regression).

Your plot should include **four lines** plotted against polynomial degree:

1. **Bias squared**
2. **Variance**
3. **Bias squared + Variance**
4. **Empirical MSE** (computed using sklearn)

Include a clear **legend** and label your axes.

This will help you visually assess how regularization impacts bias, variance, and overall model error.

(10 points)

#### A.2.5 e) Evaluate the Regularized Model

- What is the **optimal polynomial degree** for the Ridge Regression model, based on the **lowest empirical MSE**?  
(3 points)
- Report the corresponding values of:
  - **Bias squared**
  - **Variance**
  - **Empirical MSE**  
for that optimal Ridge model.  
(3 points)
- Compare these results to those of the optimal **Linear Regression** model from Question 1.  
Discuss how **regularization** influenced the **bias**, **variance**, and **overall prediction error (MSE)**.  
(4 points)

#### A.2.6 f) Interpreting the Impact of Regularization

- Was **regularization successful** in reducing the **total prediction error (MSE)** compared to the unregularized model?  
(2 points)
- Based on your results from **2e**, explain how **bias** and **variance** changed as a result of regularization.  
How did these changes affect the final total error?  
Support your explanation with values or observations from your analysis.  
(3 points)

## B Assignment 2

### Instructions

1. You may talk to a friend, discuss the questions and potential directions for solving them. However, you need to write your own solutions and code separately, and not as a group activity.
2. Write your code in the **Code cells** and your answers in the **Markdown cells** of the Jupyter notebook. Ensure that the solution is written neatly enough to for the graders to understand and follow.
3. Use [Quarto](#) to render the **.ipynb** file as HTML. You will need to open the command prompt, navigate to the directory containing the file, and use the command: `quarto render filename.ipynb --to html`. Submit the HTML file.
4. The assignment is worth 100 points, and is due on **Monday, 18th April 2025 at 11:59 pm**.
5. **Five points are properly formatting the assignment.** The breakdown is as follows:
  - Must be an HTML file rendered using Quarto (**1 point**). *If you have a Quarto issue, you must mention the issue & quote the error you get when rendering using Quarto in the comments section of Canvas, and submit the ipynb file.*
  - No name can be written on the assignment, nor can there be any indicator of the student's identity—e.g. printouts of the working directory should not be included in the final submission. (**1 point**)
  - There aren't excessively long outputs of extraneous information (e.g. no printouts of entire data frames without good reason, there aren't long printouts of which iteration a loop is on, there aren't long sections of commented-out code, etc.) (**1 point**)
  - Final answers to each question are written in the Markdown cells. (**1 point**)
  - There is no piece of unnecessary / redundant code, and no unnecessary / redundant text. (**1 point**)
6. The maximum possible score in the assignment is  $103+5 = 108$  out of 100.

## B.1 Optimizing KNN for Classification (71 points)

In this question, you will use `classification_data.csv`. Each row is a loan and the each column represents some financial information as follows:

- `hi_int_prncp_pd`: Indicates if a high percentage of the repayments went to interest rather than principal. **This is the classification response.**
- `out_prncp_inv`: Remaining outstanding principal for portion of total amount funded by investors
- `loan_amnt`: The listed amount of the loan applied for by the borrower. If at some point in time, the credit department reduces the loan amount, then it will be reflected in this value.
- `int_rate`: Interest Rate on the loan
- `term`: The number of payments on the loan. Values are in months and can be either 36 or 60.

As indicated above, `hi_int_prncp_pd` is the response and all the remaining columns are predictors. You will tune and train a K-Nearest Neighbors (KNN) classifier throughout this question.

### B.1.1 a) Load the Dataset (1 point)

Read the dataset into your notebook.

### B.1.2 b) Define Predictor and Response Variables (1 point)

Create the `predictor` (features) and `response` (target) variables from the dataset.

### B.1.3 c) Split the Data into Training and Test Sets (1 points)

Create the training and test datasets using the following specifications:

- Use a **75%-25% split**.
- Ensure that the **class ratio is preserved** in both training and test sets (i.e., stratify the split).
- Set `random_state=45` for reproducibility.

#### B.1.4 d) Check Class Ratios (2 points)

Print the **class distribution** (ratios) for:

- The entire dataset
- The training set
- The test set

This is to verify that the **class ratio is preserved** after splitting.

#### B.1.5 e) Scale the Dataset (2 points)

Use **StandardScaler** to scale the dataset in order to prepare it for KNN modeling.

Scaling ensures that all features contribute equally to the distance calculations used by the KNN algorithm.

#### B.1.6 f) Set Up Cross-Validation (2 points)

Before creating and tuning your model, you need to define cross-validation settings to ensure consistent and accurate evaluation across folds.

Please follow these specifications:

- Use **5 stratified folds** to preserve class distributions in each split.
- **Shuffle** the data before splitting to introduce randomness.
- Set **random\_state=14** for reproducibility.

**Note:** You must use these exact cross-validation settings throughout the rest of this question to maintain consistency.

#### B.1.7 g) Tune K for KNN Using Cross-Validation (12 points)

Tune a **KNN Classifier** using cross-validation with the following specifications:

- Use **every odd K value from 1 to 50** (inclusive).
- Keep all other model settings at their defaults.
- Use the **cross-validation settings** defined in part (f).
- Evaluate performance using the **F1 score** as the metric.

(4 points)

Then, complete the following tasks:

- Create a **plot of K values vs. cross-validation F1 scores** to visualize how K balances overfitting and underfitting. (2 points)
- Print the **best average cross-validation F1 score**. (1 points)
- Report the **K value corresponding to the best F1 score**. (1 points)
- Determine whether this is the **only K value** that results in the best F1 score. Use code to justify your answer. (2 points)
- Reflect on whether **accuracy** is a good metric for tuning the model in this case. Explain your reasoning. (2 points)

**Hint:**

In addition to reporting the best K and best F1 score, you may also want to examine the full cross-validation results to check if other K values achieved the same F1 score.

### **B.1.8 h) Optimize the Classification Threshold (4 points)**

Using the **optimal K value** you identified in part (g), optimize the classification **threshold** to maximize the cross-validation F1 score.

#### **B.1.8.1 Specifications:**

- Search across all possible threshold values using a **step size of 0.05**.
- Use the **cross-validation settings** defined in part (f).
- Evaluate performance using the **F1 score**, consistent with part (g).

#### **B.1.8.2 Tasks:**

- Visualize the **F1 score vs. different threshold values**. (2 points)
- Identify and report the **best threshold** that yields the highest F1 score. (1 points)
- Output the **best cross-validation F1 score**. (1 points)

### B.1.9 i) Evaluate the Tuning Method (2 points)

Is the method we used in parts (g) and (h) **guaranteed** to find the best combination of **K** and **threshold**, i.e., to tune the classifier to its optimal values?

(1 point)

Justify your answer.

(1 point)

### B.1.10 j) Evaluate Tuned Classifier on Test Set (3 points)

Using the **tuned KNN classifier** and the **optimal threshold** you identified, evaluate the model on the **test set**. Report the following metrics:

- F1 Score
- Accuracy
- Precision
- Recall
- AUC

### B.1.11 k) Jointly Tune K and Threshold (6 points)

Now, tune **K** and the **classification threshold simultaneously**, rather than sequentially.

- Use the same settings from the previous parts (i.e., odd K values from 1 to 50, threshold step size of 0.05, F1 score as the metric, and the same cross-validation strategy).
- Identify the **best F1 score**, along with the **K value and threshold** that produce it.

### B.1.12 l) Visualize Cross-Validation Results with a Heatmap (3 points)

Create a **heatmap** to visualize the cross-validation results in two dimensions.

- The **x-axis** should represent the **K values**.
- The **y-axis** should represent the **threshold values**.
- The color should represent the **F1 score**.

**Note:** This question only requires **one line of code**. You'll need to recall a **data visualization function** and a **data reshaping method** from 303-1.

### B.1.13 m) Compare Joint vs. Sequential Tuning Results (4 points)

- How does the **best cross-validation F1 score** from part (k) compare to the scores from parts (g) and (h)? **(1 point)**
- Did the **optimal K value** and **threshold** change when tuning them jointly? **(1 point)**
- Explain **why or why not**. Consider how tuning the two parameters together might impact the result. **(2 points)**

### B.1.14 n) Evaluate Final Tuned Model on Test Set (3 points)

Using the **tuned classifier and threshold** from part (k), evaluate the model on the **test set**. Report the following metrics:

- F1 Score
- Accuracy
- Precision
- Recall
- AUC

### B.1.15 o) Compare Tuning Strategies and Computational Cost (3 points)

Compare the tuning approach used in parts (g) & (h) (separate tuning of K and threshold) with the approach in **part (k)** (joint tuning of K and threshold) in terms of **computational cost**.

- How many **K and threshold combinations** did you evaluate in each approach? **(2 points)**
- Based on this comparison and your answer from part (l), explain the **main trade-off** involved in model tuning (e.g., between computation and performance). **(2 points)**

### B.1.16 p) Tune K Using Multiple Metrics (5 points)

GridSearchCV and `cross_val_score` are designed to optimize based on a **single metric**. In this section, you'll practice tuning hyperparameters while evaluating **multiple metrics** simultaneously using `cross_validate`.

For this imbalanced classification task, instead of optimizing the F1 score directly, we'll focus on **precision** and **recall** together.



Keep in mind that the F1 score is the **harmonic mean** of precision and recall—it balances the trade-off between the two.

Cross-validate a **KNN classifier** using the following specifications:

- Use the **same cross-validation setting** and **hyperparameter grid** as before
- Evaluate the model using **precision**, **recall**, and **f1-score**, as metrics **at the same time**.

Save the cross-validation results into a **DataFrame**, and compute the **average score for each metric**, and visualize how these metrics change with different values of  $K$ .

#### **B.1.17 q) Optimize for Recall with Precision Constraint (4 point)**

Identify the **K value** that yields the **highest recall**, while maintaining a **precision of at least 75%**.

**(3 points)**

Then, print the **average cross-validation metrics** (f1-score, precision, recall) for that  $K$  value.

**(1 point)**

#### **B.1.18 r) Tune Threshold for Maximum Recall (3 point)**

Using the **optimal K value** identified in part (q), find the **threshold** that maximizes **cross-validation Recall**, following the specifications below:

- Evaluate all possible threshold values with a **step size of 0.05**.
- Use the **cross-validation settings** from part (f).

Then: - Print the **best cross-validation recall**. - Report the **threshold value** that achieves this recall.

**Note:** This task is very similar to part (h), but it's important for the next part.

#### **B.1.19 s) Evaluate Precision-Optimized Model on Test Set (2 points)**

Using the **tuned classifier and threshold** from parts (q) and (r), evaluate the model on the **test set**. Report the following metrics:

- F1 Score
- Accuracy

- Precision
- Recall
- AUC

### B.1.20 t) Final Reflection: Comparing Tuning Strategies (3 points)

You have now tuned your KNN classifier using **three different strategies**:

1. **Sequential tuning** of K and threshold based on **F1 score** (parts g–h)
2. **Joint tuning** of K and threshold using **F1 score** (part k)
3. Tuning based on **multiple metrics**, selecting the K with the **highest recall** while maintaining **precision 75%** (parts p–r)

Reflect on the following:

- Which tuning strategy led to the **best overall performance on the test set**, based on the metrics you care about most?
- Which strategy would you choose in a real-world application, and why?
- What are the **trade-offs** between tuning for F1 score versus prioritizing precision or recall individually?

**Note:** This is an open-ended question. As long as your reasoning makes sense, you will receive full credit.

## B.2 Tuning a KNN Regressor on Bank Loan Data (32 points)

In this question, you will use `bank_loan_train_data.csv` to tune (*the model hyperparameters*) and train the model. Each row is a loan and the each column represents some financial information as follows:

- `money_made_inv`: Indicates the amount of money made by the bank on the loan. **This is the regression response.**
- `out_prncp_inv`: Remaining outstanding principal for portion of total amount funded by investors
- `loan_amnt`: The listed amount of the loan applied for by the borrower. If at some point in time, the credit department reduces the loan amount, then it will be reflected in this value.
- `int_rate`: Interest Rate on the loan

- **term**: The number of payments on the loan. Values are in months and can be either 36 or 60
- **mort\_acc**: The number of mortgage accounts
- **application\_type\_Individual**: 1 if the loan is an individual application or a joint application with two co-borrowers
- **tot\_cur\_bal**: Total current balance of all accounts
- **pub\_rec**: Number of derogatory public records

As indicated above, **money\_made\_inv** is the response and all the remaining columns are predictors. You will tune and train a K-Nearest Neighbors (KNN) regressor throughout this question.

### B.2.1 a) Split, Scale, and Tune a KNN Regressor (15 point)

Create the **training and test datasets** using the following specifications:

- Use an **80%-20% split**.
- Set **random\_state=1** for reproducibility.

Then, **scale all the predictors**, as KNN is sensitive to the scale of input features.

Next, you will **tune a KNN Regressor** by searching for the optimal hyperparameters using three search approaches: **Grid Search**, **Random Search**, and **Bayesian Search**.

#### B.2.1.1 Cross-Validation Setting

You should use **5-fold cross-validation**, with the following specifications:

- The data should be **shuffled** before splitting
- Use **random\_state=1** to ensure **reproducibility**

### B.2.1.2 Hyperparameters to Tune:

You will tune the following hyperparameters for the KNN Regressor:

You will tune the following hyperparameters for the **K-Nearest Neighbors Regressor**, using **Minkowski** as the distance metric:

1. **n\_neighbors**: Number of nearest neighbors
  - Tune over the range: `np.arange(1, 25, 1)`
2. **p**: Power parameter for the Minkowski distance
  - Use values: `np.arange(1, 4, 1)`
  - `p = 1` corresponds to Manhattan distance
  - `p = 2` corresponds to Euclidean distance
  - **Note**: Set the distance metric to "minkowski"
3. **weights**: Weight function used in prediction  
You must consider the following **5 types of weights**:
  - 'uniform': All neighbors are weighted equally
  - 'distance': Weight is inversely proportional to distance
  - Custom weight functions:
    - $\propto \frac{1}{\text{distance}^2}$
    - $\propto \frac{1}{\text{distance}^3}$
    - $\propto \frac{1}{\text{distance}^4}$

For **each search method** (Grid Search, Random Search, Bayesian Search), report the following:

- **best\_params\_**: The best combination of hyperparameters
- **best\_score\_**: Cross-validated RMSE on the training set
- **Test RMSE** obtained from the best model
- **Execution time** for the search process

**Hint:**

Define **three custom weight functions** as shown below:

```
def dist_power_2(distance):
    return 1 / (1e-10 + distance**2)

def dist_power_3(distance):
    return 1 / (1e-10 + distance**3)

def dist_power_4(distance):
    return 1 / (1e-10 + distance**4)
```

Note the small constant `1e-10` helps avoid division by zero and numerical instability.

### B.2.2 b) Compare Tuning Approaches (1 point)

Compare the results from part (2a) in terms of **execution time** and **model performance**. Briefly discuss the **main trade-offs** among the three hyperparameter tuning approaches: Grid Search, Random Search, and Bayesian Search.

### B.2.3 c) Feature Selection and Hyperparameter Tuning with GridSearchCV (15 point)

KNN performance can **deteriorate significantly** if irrelevant or noisy predictors are included. In this part, you will explore **feature selection** to improve model performance, followed by **hyperparameter tuning** using GridSearchCV (with `refit=True`).

Try the following **three different feature selection approaches**:

1. **Correlation-based filtering:**

- Select features with an absolute correlation of at least **0.1** with the target variable.

2. **Lasso regression for feature selection:**

- Use `Lasso(alpha=50)` to select important features based on non-zero coefficients.

3. **SelectKBest:**

- Use `SelectKBest` with `f_regression`, selecting the **top 4** features.

For **each approach**, perform hyperparameter tuning using **GridSearchCV**, and report:

- The **best score** (cross-validated RMSE) on the **training set**
- The **test RMSE** from the best model
- The **best hyperparameters**

#### B.2.4 d) Compare Feature Selection Approaches (1 point)

Create a **DataFrame** that summarizes the model performance from each feature selection method, including:

- **Training RMSE**
- **Test RMSE**

Be sure to also include the results from the model trained **without any feature selection** for comparison.

Then, briefly explain what you learned from this experiment.

For example: Did feature selection improve performance? Which method worked best?

# C Assignment 3

## Instructions

1. You may talk to a friend, discuss the questions and potential directions for solving them. However, you need to write your own solutions and code separately, and not as a group activity.
2. Write your code in the *Code* cells and your answer in the *Markdown* cells of the Jupyter notebook. Ensure that the solution is written neatly enough to understand and grade.
3. Use [Quarto](#) to print the *.ipynb* file as HTML. You will need to open the command prompt, navigate to the directory containing the file, and use the command: `quarto render filename.ipynb --to html`. Submit the HTML file.
4. The assignment is worth 100 points, and is due on **Friday, 2th May 2025 at 11:59 pm**.
5. **Five points are properly formatting the assignment.** The breakdown is as follows:
  - Must be an HTML file rendered using Quarto (2 pts). *If you have a Quarto issue, you must mention the issue & quote the error you get when rendering using Quarto in the comments section of Canvas, and submit the ipynb file. If your issue doesn't seem genuine, you will lose points.*
  - There aren't excessively long outputs of extraneous information (e.g. no printouts of entire data frames without good reason, there aren't long printouts of which iteration a loop is on, there aren't long sections of commented-out code, etc.) (1 pt)
  - Final answers of each question are written in Markdown cells (1 pt).
  - There is no piece of unnecessary / redundant code, and no unnecessary / redundant text (1 pt)

## C.1 Regression Problem - Miami housing

### C.1.1 a) Data preparation

Read the data *miami-housing.csv*. Check the description of the variables [here](#). Split the data into 60% train and 40% test. Use `random_state = 45`. The response is `SALE_PRC`, and the rest

of the columns are predictors, except `PARCELNO`. Print the shape of the predictors dataframe of the train data.

*(2 points)*

### C.1.2 b) Baseline Decision Tree Model

Train a **Decision Tree Regressor** to predict `SALE_PRC` using all available predictors.

- Use `random_state=45` and keep all other hyperparameters at their default values.
- After training the model, evaluate and report the following on **both the training and test sets**:
  - **Mean Absolute Error (MAE)**
  - **R<sup>2</sup> Score**

*(3 points)*

### C.1.3 c) Tune the Decision Tree Model

Tune the hyperparameters of the **Decision Tree Regressor** developed in the previous question and evaluate its performance.

Your goal is to achieve a **test set MAE (Mean Absolute Error) below \$68,000**.

- You must display the **optimal hyperparameter values** obtained from the tuning process.
- Compute and report the **test MAE and R<sup>2</sup> Score** using the tuned model.

**Hints:** 1. `BayesSearchCV()` with `max_depth` and `max_features` can often complete in under a minute. 2. You may use **any hyperparameter tuning method** (e.g., `GridSearchCV`, `RandomizedSearchCV`, `BayesSearchCV`). 3. You are free to choose **which hyperparameters to tune** and define your own **search space**.

*(9 points)*

### C.1.4 d) Bagged Decision Trees with Out-of-Bag Evaluation

Train a **Bagging Regressor** using Decision Trees as base estimators to predict `SALE_PRC`.

- Use only the `n_estimators` hyperparameter for tuning; keep all other parameters at their default values.
- Increase the number of trees (`n_estimators`) until the **out-of-bag (OOB) MAE stabilizes**.



- Report the final **OOB MAE**, **test MAE**, and **R<sup>2</sup> Score**, and ensure that the OOB MAE is less than \$48,000.

*(4 points)*

### C.1.5 e) Bagged Decision Trees Without Bootstrapping

Train a **Bagging Regressor** using Decision Trees, but this time **disable bootstrapping** by setting `bootstrap=False`.

- Use the same `n_estimators` value as in the previous question.
- Compute and report the following on the **test set**:
  - **Mean Absolute Error (MAE)**
  - **R<sup>2</sup> Score**

Explain **why the test MAE** in this case is:

- **Much higher** than the MAE obtained when bootstrapping was enabled (previous question).
- **Lower** than the MAE obtained from a single untuned decision tree (as in Question 1(b)).

Hint: Consider the impact of bootstrap sampling on variance reduction and the benefits of aggregation in ensemble methods.

*(2 point for code, 3 + 3 points for reasoning)*

### C.1.6 f) Bagged Decision Trees with Feature Bootstrapping Only

Train a **Bagging Regressor** using Decision Trees, with the following configuration: - **Disable sample bootstrapping** by setting `bootstrap=False` - **Enable feature bootstrapping** by setting `bootstrap_features=True`

Use the same number of estimators (`n_estimators`) as in the previous bagging experiments.

- Compute and report the following on the **test set**:
  - **Mean Absolute Error (MAE)**
  - **R<sup>2</sup> Score**

Explain why the **test MAE** obtained in this setting is **much lower** than the one in the previous question, where neither bootstrapping samples nor features was used.

*(2 point for code, 3 points for reasoning)*

### C.1.7 g) Tuning a Bagged Tree Model

#### C.1.7.1 i) Approaches

There are two common approaches for tuning a **bagged tree model**:

1. **Out-of-Bag (OOB) Prediction**
2.  **$K$ -fold Cross-Validation** using GridSearchCV

What is the advantage of each approach over the other? Specifically:

- What is the **advantage of the out-of-bag approach** compared to  $K$ -fold cross-validation?
- What is the **advantage of  $K$ -fold cross-validation** compared to the out-of-bag approach?

*(3 + 3 points)*

#### C.1.7.2 ii) Tuning the hyperparameters

Tune the hyperparameters of the bagged tree model developed in 1(d). You may use either of the approaches mentioned in the previous question. Show the optimal values of the hyperparameters obtained. Compute the MAE and  $R^2$  Score on test data with the tuned model. **Your test MAE must be less than the test MAE obtained in the previous question.**

It is up to you to pick the hyperparameters and their values in the grid.

**Hint:**

GridSearchCV() may work better than BayesSearchCV() in this case.

*(9 points)*

### C.1.8 h) Random Forest

#### C.1.8.1 i) Tuning a Random Forest Model

Train and tune a **Random Forest Regressor** to predict SALE\_PRC.

- Select hyperparameters and define your own tuning grid.
- Use any tuning approach (e.g., Out-of-Bag (OOB) evaluation or  $K$ -fold cross-validation).
- Report the following performance metrics on the **test set**:
  - **Mean Absolute Error (MAE)**
  - **$R^2$  Score**

Your goal is to achieve a **test MAE below \$46,000**.

**Hint:**

The **OOB approach** is efficient and can complete in under a minute.

*(9 points)*

### C.1.8.2 ii) Feature Importance

After fitting the tuned **Random Forest Regressor**, extract and display the **feature importances**.

- Print the predictors in **decreasing order of importance** based on the trained model.
- This helps identify which features contribute most to predicting **SALE\_PRC**.

*(4 points)*

### C.1.8.3 iii) Feature Selection

Drop the **least important predictor** identified in the previous step, and re-train the **tuned Random Forest model**.

- Compute the **test MAE and R<sup>2</sup> Score** after dropping the feature.
- You may need to adjust the **max\_features** hyperparameter to reflect the reduced number of predictors.
- Compare the new test MAE with the previous one.

Did the test MAE decrease after removing the least important feature?

*(4 points)*

### C.1.8.4 iv) Random Forest vs. Bagging: max\_features

The **max\_features** hyperparameter is available in both **RandomForestRegressor()** and **BaggingRegressor()**.

Does **max\_features** have the **same meaning** in both models?

If not, explain the **difference in how it is interpreted and applied**.

**Hint:** Refer to the scikit-learn documentation for both estimators to understand how **max\_features** affects feature selection during training.

*(1 + 3 points)*

## C.2 Classification - Term deposit

The data for this question is related with direct marketing campaigns of a Portuguese banking institution. The marketing campaigns were based on phone calls, where bank clients were called to subscribe for a term deposit.

There is a train data - *train.csv*, which you will use to develop a model. There is a test data - *test.csv*, which you will use to test your model. Each dataset has the following attributes about the clients called in the marketing campaign:

1. **age**: Age of the client
2. **education**: Education level of the client
3. **day**: Day of the month the call is made
4. **month**: Month of the call
5. **y**: did the client subscribe to a term deposit?
6. **duration**: Call duration, in seconds. This attribute highly affects the output target (e.g., if **duration**=0 then **y**='no'). Yet, the duration is not known before a call is performed. Also, after the end of the call **y** is obviously known. Thus, this input should only be included for inference purposes and should be discarded if the intention is to have a realistic predictive model.

(Raw data source: [Source](#). Do not use the raw data source for this assignment. It is just for reference.)

### C.2.1 a) Data Preparation

Begin by examining the **distribution of the target variable** in both the training and test sets. This will help you assess whether there is any significant **class imbalance**.

Next, consider the two available approaches for hyperparameter tuning:

- **Cross-validation (CV)**
- **Out-of-bag (OOB) evaluation**

### C.2.1.1 Which method do you prefer for this dataset, and why?

Discuss your choice based on:

- The **size of the dataset**
- The **class imbalance** in the target variable
- The **reliability and interpretability** of each method
- Whether you need **stratified sampling** to preserve class distribution during evaluation

(2 points)

### C.2.2 b) Random Forest for Term Deposit Subscription Prediction

Develop and tune a **Random Forest Classifier** to predict whether a client will subscribe to a term deposit using the following predictors:

- `age`
- `education`
- `day`
- `month`

The model must satisfy the following performance criteria:

#### C.2.2.1 Requirements:

1. **Minimum overall classification accuracy of 75%**, across both *train.csv* and *test.csv*.
2. **Minimum recall of 60%**, across both *train.csv* and *test.csv*.

You must:

- Print the **accuracy** and **recall** for both datasets (*train.csv* and *test.csv*).
- Use **cross-validation on the training data** to optimize the model hyperparameters.
- Select a **threshold probability** for classification and apply it consistently across both datasets.

#### C.2.2.2 Important Notes:

- i. **Do not use duration** as a predictor. Its value is determined after the marketing call ends, so using it would leak information about the outcome.
- ii. You are free to choose any **decision threshold** for classification, but the same threshold must be used consistently for both training and test evaluation.

- iii. Use **cross-validation** to tune hyperparameters such as `max_features`, `max_depth`, and `max_leaf_nodes`.
  - You may use `StratifiedKFold` or any appropriate CV method that respects class imbalance.
- iv. After tuning the model, **plot cross-validated accuracy and recall** across a range of threshold values (e.g., 0.1 to 0.9). Use this plot to select a threshold that satisfies the required trade-off between accuracy and recall.
- v. **Evaluate the final tuned model (with the chosen threshold)** on the test dataset. Do not use the test data to guide any part of the tuning or threshold selection.

### C.2.2.3 Hints:

- Restrict the search space to:
  - `max_depth` 25
  - `max_leaf_nodes` 45These limits encourage generalization and help balance recall and accuracy.
- Consider using cross-validation scores to compute predicted probabilities when plotting recall/accuracy curves.

### C.2.2.4 Scoring Breakdown (22 points total):

- **8 points** – Hyperparameter tuning via cross-validation
- **5 points** – Plotting accuracy and recall across thresholds
- **5 points** – Threshold selection based on the plot
- **4 points** – Reporting accuracy and recall on both datasets

## C.3 Predictor Transformations in Trees

Can a **non-linear monotonic transformation** of predictors (such as `log()`, `sqrt()`, etc.) be useful in improving the accuracy of **decision tree models**?

Provide a brief explanation based on your understanding of how decision trees split data and handle predictor scales.

*(4 points for answer)*

## D Datasets, assignment and project files

Datasets used in the book, assignment files, project files, and prediction problems report template can be found [here](#)