

Random_Forest_Classifiers_Weber_Constant

October 8, 2023

1 Temporal Metrics: Quantifying Human Time Perception

Regis University Cammie Newmyer October 8, 2023

2 Purpose Statement: Quantifying Human Time Perception

Objective: To introduce an innovative approach to understanding and quantifying human time perception, bridging the realms of neurology, physics, and psychology to offer new insights into a long-standing question: How do humans perceive the passage of time?

The Importance of Understanding Human Time Perception: The way humans perceive time has far-reaching implications not only for individual subjective experiences but also for broader social and functional contexts. Our perception of time shapes our reactions, decisions, memories, and future anticipations. It influences our emotional state, the rhythm of our daily activities, and even our cultural narratives. From waiting for a bus to recollecting past experiences, our sense of time pervades every aspect of our lives. Thus, understanding it is pivotal for both enhancing personal well-being and addressing societal challenges.

A Novel Framework: The Distance = Rate x Time Paradigm in Time Perception: At first glance, the formula “distance = rate x time” seems exclusive to the physical realm, typically associated with motion. However, we postulate that it can be metaphorically applied to human time perception. Here, ‘distance’ doesn’t signify a physical journey but represents the cognitive journey our brains undertake in a 24-hour span, as they process myriad stimuli and experiences.

Distance: The theoretical “distance” our brains travel is a cumulative representation of all brainwave activities in a typical day. This cognitive journey can vary significantly among individuals, especially when considering conditions or disorders that impact time perception.

Rate: This is characterized by the wavelengths and frequencies of standard brainwaves - delta, theta, alpha, low beta, mid beta, high beta, and gamma. Each type of wave represents different cognitive states and functions, from deep sleep to heightened alertness.

Time: Time, in this context, refers to the duration one experiences a particular brainwave activity within 24 hours.

The Weber-Fechner law, or Weber’s law: A principle in experimental psychology proposed by Ernst Heinrich Weber in 1834. It states that the just-noticeable difference between two stimuli is proportional to the magnitude of the stimuli. That is, if you increase a stimulus (like brightness, weight, loudness, etc.), the amount of change required for us to notice this change becomes larger as the stimulus itself becomes larger. In the context of time perception, applying Weber’s law would

mean that our ability to distinguish between two durations would be based on a ratio or proportion rather than a fixed quantity.

Incorporating Weber’s Time Constant and Myelination: The Weber time constant, derived by calculating the difference between the ‘distance’ covered by an average brain compared to those with altered time perception, offers a scalar measure of time perception variations. Furthermore, factoring in changes in myelination, which can influence the speed and efficiency of neural transmissions, adds depth to our understanding. An increase in myelination can speed up neural processes, possibly leading to altered time perception, while a decrease might have the opposite effect.

3 Data Description

Methodology: Harnessing AI-Powered Insights for Understanding Human Time Perception

Data Acquisition:

Theoretical Judgements: Our research began by tapping into the theoretical frameworks of neuroscience, psychology, and time perception. By leveraging these well-established theories, we constructed a foundational understanding upon which more specific and nuanced data points could be built.

Rapid Fire Questioning with AI ChatGPT-4: To delve deeper into the complexities of time perception, we engaged in extensive interactions with OpenAI’s ChatGPT-4. This state-of-the-art AI model, trained on vast amounts of data, provided us with nuanced insights and knowledge gaps in existing research.

Why ChatGPT-4?

Access to Comprehensive Data: ChatGPT-4 has been trained on a multitude of research papers, articles, and databases. Its knowledge spans a wide array of fields, allowing us to extract valuable information on brainwave activity, associated tasks, and conditions influencing time perception.

Efficient Interaction: Rapid-fire questioning with ChatGPT-4 ensured that our data acquisition process was both thorough and efficient. By posing sequential questions and building upon the AI’s responses, we were able to derive detailed and interconnected insights within a short time frame.

Dynamic Learning Approach: ChatGPT-4’s ability to understand and respond contextually enabled a more organic, conversational approach to data gathering. This dynamic interaction often led to the revelation of unexpected but valuable insights.

Data Analysis:

Brainwave Activity and Associated Tasks: With the data acquired, we charted out a comprehensive mapping of different brainwave types (delta, theta, alpha, beta, gamma) against their associated cognitive tasks. This provided a clear picture of how various activities or states of being could influence time perception.

Alterations in Human Time Perception: Further, we analyzed the factors leading to alterations in time perception. Data indicated a range of influences from biological (e.g., neurochemical changes, myelination variations) to psychological (e.g., trauma, mindfulness practices).

Conclusion:

By adopting this multifaceted approach, we aim to shed light on the intricate mosaic of human time perception. Unveiling its mysteries could lead to therapeutic breakthroughs for disorders affecting time perception and offer everyone deeper insights into their own experiences of the world. Our methodology, which combined theoretical judgments with advanced AI-powered interactions, led to a rich dataset on human time perception. By leveraging the power of ChatGPT-4 and its extensive training, we have gathered insights that push the boundaries of existing knowledge and pave the way for future research in this intriguing field.

4 Feature Descriptions

Analyzed Conditions: ADHD, Aging 18-29, Aging 30-49, Aging 50-69, Aging 70-89, Aging 90+, Alcohol, Alzheimer’s Disease, Anxiety, Autism Type 1, Autism Type 2, Average, Bipolar Depressive, Bipolar Manic, Brain Damage, Brain Lesions, Caffeine, Chronic Pail disorder, Cocaine, Concentration Meditation, Depression (MDD), Dissociative Disorders, Elite Athlete, Emregency Event, Epilepsy, Fentanyl, Flow State, Graduate Student, Heroin, High IQ, Insomnia, Ketamione, Learning Problems, Low IQ, LSD, Marijuana, MDMA (Ecstasy), Methamphetamine, Migraine, Morphine, MS, Musician, Nicotine, Oxycodone, Parkinson’s Disease, Percrption Antidepressants, Percrption Sleep Aids, Percrption Stimulants, Psilocybin, Psychosis, PTSD, Savant Type 1, Savant Type 2, Schizophrenia, Severe Cognitive Disability, Stress, TBI, Tibetyan Meditation, Transcendental Meditation. Condition Features: Time Feels Faster or Slower (0 Average, 1 Both, 2 Faster, 3 Slower); Myelin Increase or Decrease (100% is Average); Speed (average speed associated with adequate myelination (75 meters per second); Speed with Myelination Increase or Decrease applied; Wavelength in meters = velocity divided by frequency; Number of Cycles = frequency times time; Total Distance (Delta waves) = wavelength times the number of cycles in kilometers; Delta Time (time in hours in an average day that a person with a given condition spends in that brainwave activity); Delta Frequency (the average frequency of delta waves); Items 4-9 repeat for theta waves, alpha waves, low beta waves, mid beta waves, high beta waves, and gamma waves; Total Time (24 hours); Cumulative Distance in kilometers; Calculated percent increase/decrease based on average person total distance; Rate = kilometers divided by hours; Weber’s Constant (altered perception of time).

Feature names:

```
feature_names = 'Target','Myelin', 'Speed', 'Speed_M', 'D_WL', 'D_Cycle', 'D_Dist',
'D_Time', 'D_Freq', 'T_WL', 'T_Cycle', 'T_Dist', 'T_Time', 'T_Freq', 'A_WL', 'A_Cycle',
'A_Dist', 'A_Time', 'A_Freq', 'BL_WL', 'BL_Cycle', 'BL_Dist', 'BL_Time', 'BL_Freq',
'BM_WL', 'BM_Cycle', 'BM_Dist', 'BM_Time', 'BM_Freq', 'BH_WL', 'BH_Cycle', 'BH_Dist',
'BH_Time', 'BH_Freq', 'G_WL', 'G_Freq', 'G_Dist', 'G_Time', 'G_Freq.1', 'Total_Time', 'To-
tal_Distance', 'Percent_Increase', 'Total_Rate', 'Weber_k'
```

Cumulative Time and Initial Speed features are dropped becuase they are constants.

5 Machine Learning Approach

The utilization of Random Forest models in this research project stems from their inherent advantages in addressing the complexities of our multiclass classification problem. Random Forest models offer a robust and insightful approach to gaining a deeper understanding of the factors influencing classification accuracy. Their applicability was chosen based on the recognition that

they provide a unique opportunity to unravel the importance of specific features, thereby informing subsequent feature selection and engineering efforts for the development of a deep learning neural network (NN).

Another compelling reason for embracing Random Forest models is their transparency in interpreting complex datasets. This transparency facilitates a clear comprehension of feature interactions and their contributions to overall accuracy. Beyond interpretability, these models serve as effective benchmarking tools, allowing for the establishment of performance baselines that aid in assessing the value of more intricate deep learning models. Moreover, their adeptness at handling class imbalances, often encountered in multiclass classification tasks, adds to their appeal, equipping us with the expertise needed to address real-world data scenarios effectively.

In summary, the choice to employ Random Forest models is a strategic one, driven by their capacity to provide critical insights into feature importance, enhance data interpretability, support model benchmarking, and tackle class imbalances. These models lay the groundwork for informing the design of a deep learning NN model capable of excelling in real-world scenarios with larger and more representative datasets.

6 Import packages and load data

```
[1]: #general
import pandas as pd
import matplotlib
import matplotlib.pyplot as plt
import numpy as np
import random as rnd

# visualization
import seaborn as sns
import matplotlib.pyplot as plt
import graphviz
%matplotlib inline
sns.set()

# sklearn packages
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import MinMaxScaler
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC, LinearSVC
from sklearn.ensemble import RandomForestClassifier, VotingClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.linear_model import Perceptron
from sklearn.linear_model import SGDClassifier
```

```

from sklearn.tree import DecisionTreeClassifier, export_graphviz
from sklearn.model_selection import GridSearchCV, cross_val_score
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import GradientBoostingClassifier

# imbalanced-learn
from imblearn.over_sampling import SMOTE

import warnings
warnings.filterwarnings("ignore")

```

```

[2]: # Load the data
df = pd.read_csv('Temporal_Metrics.csv')

df['Condition'] = df['Condition'].astype('str')

# Initial data exploration
print(df.head())
print(df.info())
print(df.describe())

# Check for missing values
print(df.isnull().sum())

```

	Condition	Target	Myelin	Speed	Speed_M	D_WL	D_Cycle	D_Dist	D_Time	\
0	ADHD	1	0.90	75.0	67.50	25.0	43200.0	1080.0	4.0	
1	Aging 18-29	0	1.00	75.0	75.00	25.0	43200.0	1080.0	4.0	
2	Aging 30-49	0	0.98	75.0	73.50	25.0	48600.0	1215.0	4.5	
3	Aging 50-69	3	0.94	75.0	70.50	25.0	54000.0	1350.0	5.0	
4	Aging 70-89	3	0.85	75.0	63.75	30.0	40500.0	1215.0	4.5	

	D_Freq	...	G_WL	G_Freq	G_Dist	G_Time	G_Freq.1	Total_Time	\
0	3.0	...	1.69	144000.0	243.0	1.0	40.0	24.0	
1	3.0	...	1.88	144000.0	270.0	1.0	40.0	24.0	
2	3.0	...	2.10	63000.0	132.3	0.5	35.0	24.0	
3	3.0	...	2.35	54000.0	126.9	0.5	30.0	24.0	
4	2.5	...	2.28	0.0	0.0	0.0	28.0	24.0	

	Total_Distance	Percent_Increase	Total_Rate	Weber_k
0	6169.50	-0.05	257.06	-1.21
1	6480.00	0.00	270.00	0.00
2	6417.90	-0.01	267.41	-0.23
3	6293.70	-0.03	262.24	-0.71
4	6014.25	-0.08	250.59	-1.86

```

[5 rows x 45 columns]
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 59 entries, 0 to 58

```

Data columns (total 45 columns):

#	Column	Non-Null Count	Dtype
0	Condition	59 non-null	object
1	Target	59 non-null	int64
2	Myelin	59 non-null	float64
3	Speed	59 non-null	float64
4	Speed_M	59 non-null	float64
5	D_WL	59 non-null	float64
6	D_Cycle	59 non-null	float64
7	D_Dist	59 non-null	float64
8	D_Time	59 non-null	float64
9	D_Freq	59 non-null	float64
10	T_WL	59 non-null	float64
11	T_Cycle	59 non-null	float64
12	T_Dist	59 non-null	float64
13	T_Time	59 non-null	float64
14	T_Freq	59 non-null	float64
15	A_WL	59 non-null	float64
16	A_Cycle	59 non-null	float64
17	A_Dist	59 non-null	float64
18	A_Time	59 non-null	float64
19	A_Freq	59 non-null	float64
20	BL_WL	59 non-null	float64
21	BL_Cycle	59 non-null	float64
22	BL_Dist	59 non-null	float64
23	BL_Time	59 non-null	float64
24	BL_Freq	59 non-null	float64
25	BM_WL	59 non-null	float64
26	BM_Cycle	59 non-null	float64
27	BM_Dist	59 non-null	float64
28	BM_Time	59 non-null	float64
29	BM_Freq	59 non-null	float64
30	BH_WL	59 non-null	float64
31	BH_Cycle	59 non-null	float64
32	BH_Dist	59 non-null	float64
33	BH_Time	59 non-null	float64
34	BH_Freq	59 non-null	float64
35	G_WL	59 non-null	float64
36	G_Freq	59 non-null	float64
37	G_Dist	59 non-null	float64
38	G_Time	59 non-null	float64
39	G_Freq.1	59 non-null	float64
40	Total_Time	59 non-null	float64
41	Total_Distance	59 non-null	float64
42	Percent_Increase	59 non-null	float64
43	Total_Rate	59 non-null	float64
44	Weber_k	59 non-null	float64

dtypes: float64(43), int64(1), object(1)

memory usage: 20.9+ KB

None

	Target	Myelin	Speed	Speed_M	D_WL	D_Cycle \
count	59.000000	59.000000	59.0	59.000000	59.000000	59.000000
mean	2.067797	0.914746	75.0	68.605932	25.514746	45289.830508
std	1.048224	0.112593	0.0	8.444444	3.456182	16243.425141
min	0.000000	0.500000	75.0	37.500000	21.430000	10800.000000
25%	1.000000	0.880000	75.0	66.000000	25.000000	37800.000000
50%	2.000000	0.900000	75.0	67.500000	25.000000	43200.000000
75%	3.000000	0.990000	75.0	74.250000	25.000000	54000.000000
max	3.000000	1.100000	75.0	82.500000	37.500000	88200.000000

	D_Dist	D_Time	D_Freq	T_WL	...	G_WL \
count	59.000000	59.000000	59.000000	59.000000	...	59.000000
mean	1116.610169	4.135593	2.983051	12.307458	...	1.786441
std	323.796766	1.199247	0.334330	1.177927	...	0.320537
min	405.000000	1.500000	2.000000	10.710000	...	1.070000
25%	1012.500000	3.750000	3.000000	11.540000	...	1.590000
50%	1080.000000	4.000000	3.000000	12.500000	...	1.780000
75%	1350.000000	5.000000	3.000000	12.500000	...	1.990000
max	1890.000000	7.000000	3.500000	16.670000	...	2.360000

	G_Freq	G_Dist	G_Time	G_Freq.1	Total_Time \
count	59.000000	59.000000	59.000000	59.000000	59.0
mean	105955.932203	175.591525	0.703390	39.457627	24.0
std	67695.199930	90.855849	0.336291	7.728858	0.0
min	0.000000	0.000000	0.000000	25.000000	24.0
25%	63000.000000	121.500000	0.500000	35.000000	24.0
50%	72000.000000	132.300000	0.500000	40.000000	24.0
75%	144000.000000	243.000000	1.000000	41.000000	24.0
max	360000.000000	513.000000	2.000000	60.000000	24.0

	Total_Distance	Percent_Increase	Total_Rate	Weber_k
count	59.000000	59.000000	59.000000	59.000000
mean	6230.684746	-0.042881	259.612373	-1.031017
std	324.993176	0.057237	13.541432	1.375004
min	5238.000000	-0.240000	218.250000	-5.690000
25%	6102.000000	-0.060000	254.250000	-1.490000
50%	6234.300000	-0.040000	259.760000	-0.950000
75%	6448.950000	-0.005000	268.705000	-0.115000
max	6817.500000	0.050000	284.060000	1.190000

[8 rows x 44 columns]

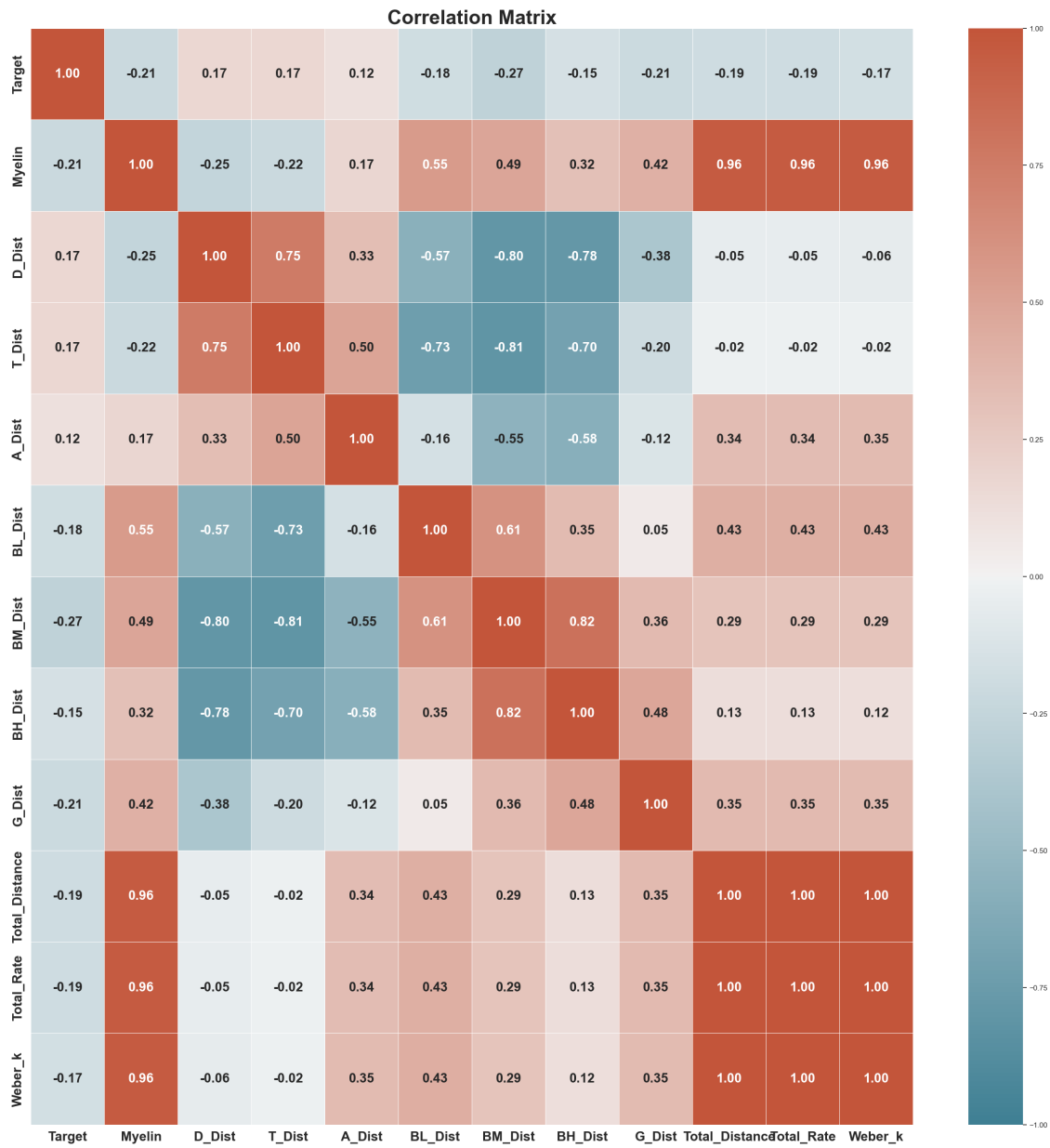
Condition	0
Target	0
Myelin	0
Speed	0

Speed_M	0
D_WL	0
D_Cycle	0
D_Dist	0
D_Time	0
D_Freq	0
T_WL	0
T_Cycle	0
T_Dist	0
T_Time	0
T_Freq	0
A_WL	0
A_Cycle	0
A_Dist	0
A_Time	0
A_Freq	0
BL_WL	0
BL_Cycle	0
BL_Dist	0
BL_Time	0
BL_Freq	0
BM_WL	0
BM_Cycle	0
BM_Dist	0
BM_Time	0
BM_Freq	0
BH_WL	0
BH_Cycle	0
BH_Dist	0
BH_Time	0
BH_Freq	0
G_WL	0
G_Freq	0
G_Dist	0
G_Time	0
G_Freq.1	0
Total_Time	0
Total_Distance	0
Percent_Increase	0
Total_Rate	0
Weber_k	0
dtype: int64	

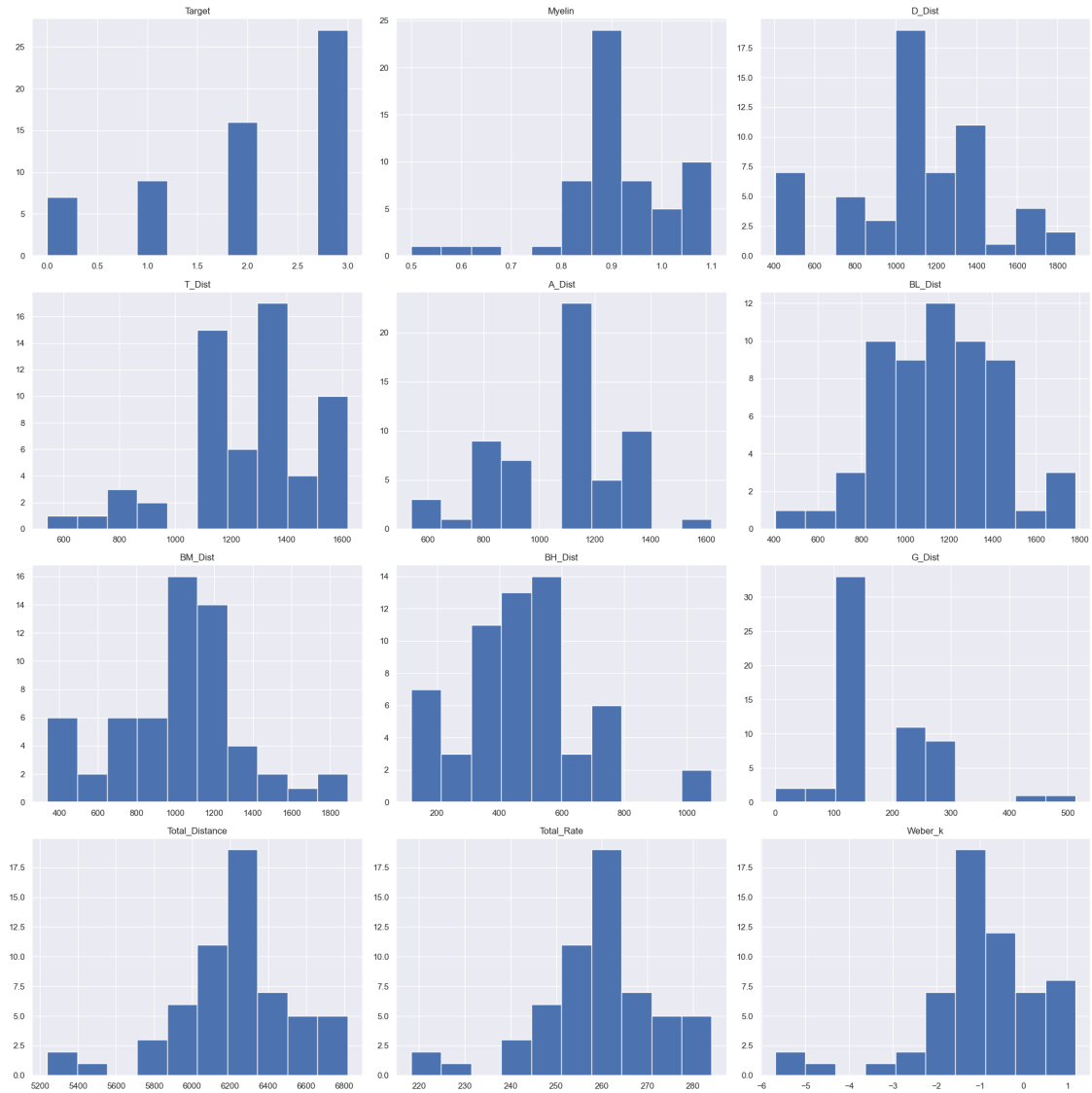
7 Exploratory Data Analysis

```
[3]: # Correlation matrix
corr_matrix = df[['Target', 'Myelin', 'D_Dist',
    ↳ 'T_Dist', 'A_Dist', 'BL_Dist', 'BM_Dist', 'BH_Dist', 'G_Dist', 'Total_Distance', 'Total_Rate', 'Web
    ↳ corr(numeric_only=True)

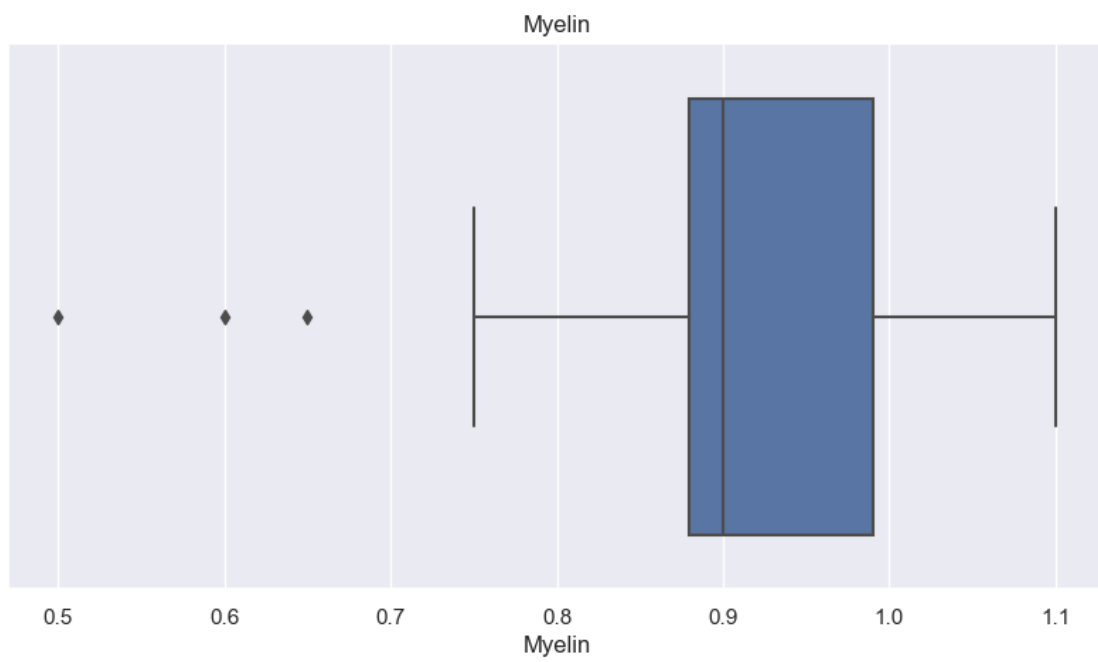
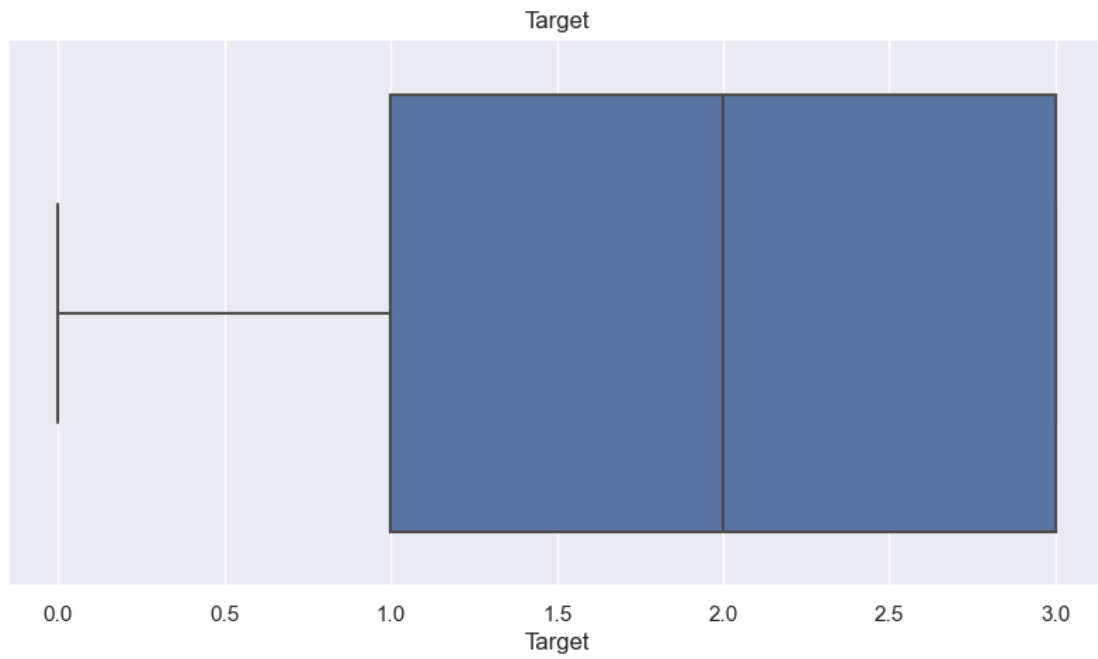
# Plotting
plt.figure(figsize=(30,30))
sns.heatmap(
    corr_matrix,
    annot=True,
    fmt=".2f",
    cmap=sns.diverging_palette(220, 20, as_cmap=True),
    vmin=-1,
    vmax=1,
    annot_kws={"size": 20, "weight": "bold"}, # Adjust font size and weight
    ↳ for the annotations
    linewidths=0.5 # Adjusts line width for gridlines
)
plt.xticks(fontsize=20, weight='bold') # Adjust x-tick label font size and
    ↳ weight
plt.yticks(fontsize=20, weight='bold') # Adjust y-tick label font size and
    ↳ weight
plt.title('Correlation Matrix', fontsize=30, weight='bold')
plt.show()
```

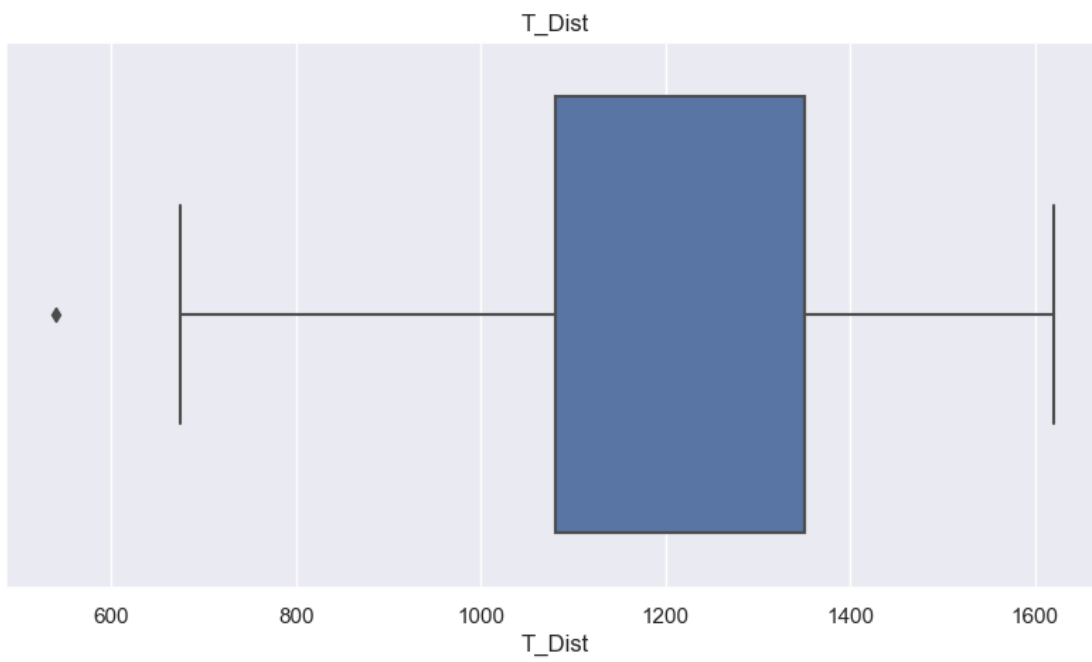
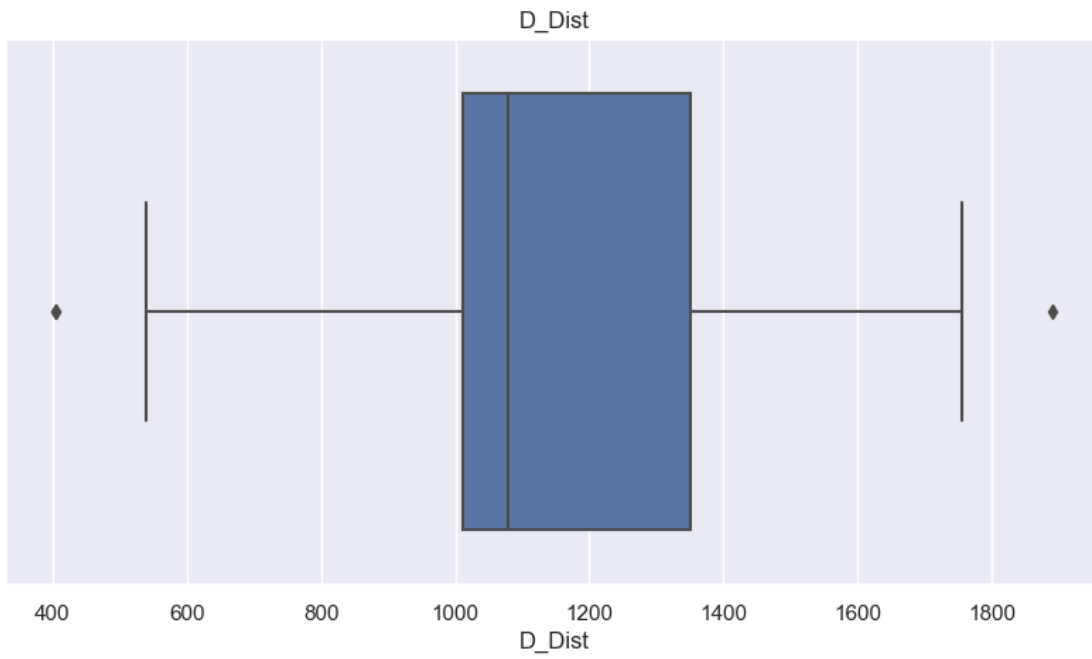


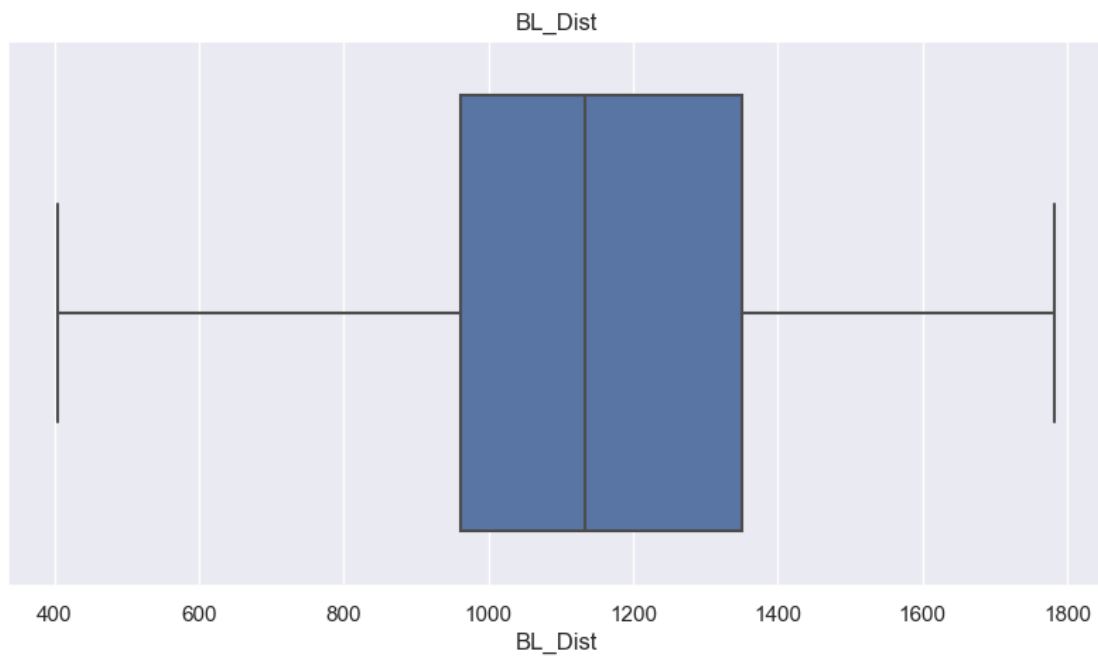
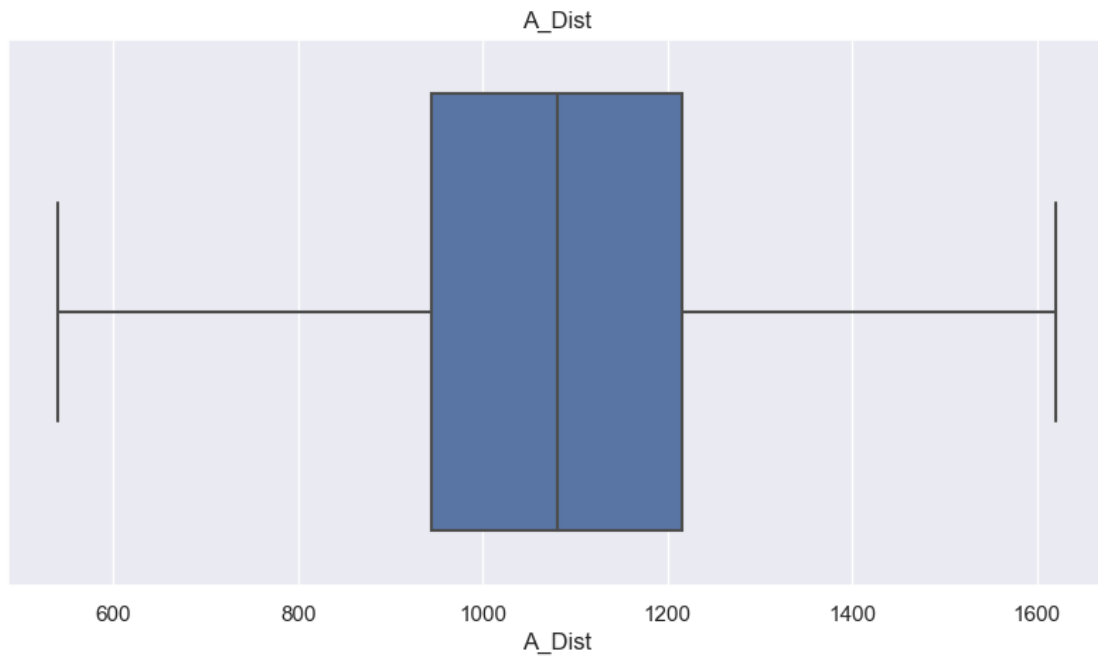
```
[4]: # Visualization
# Histogram for all features
df[['Target', 'Myelin', 'D_Dist', 'T_Dist', 'A_Dist', 'BL_Dist', 'BM_Dist', 'BH_Dist', 'G_Dist', 'Total_Distance', 'Total_Rate', 'Weber_k']]
hist(figsize=(20,20))
plt.tight_layout() # Ensures that plots don't overlap
plt.show()
```

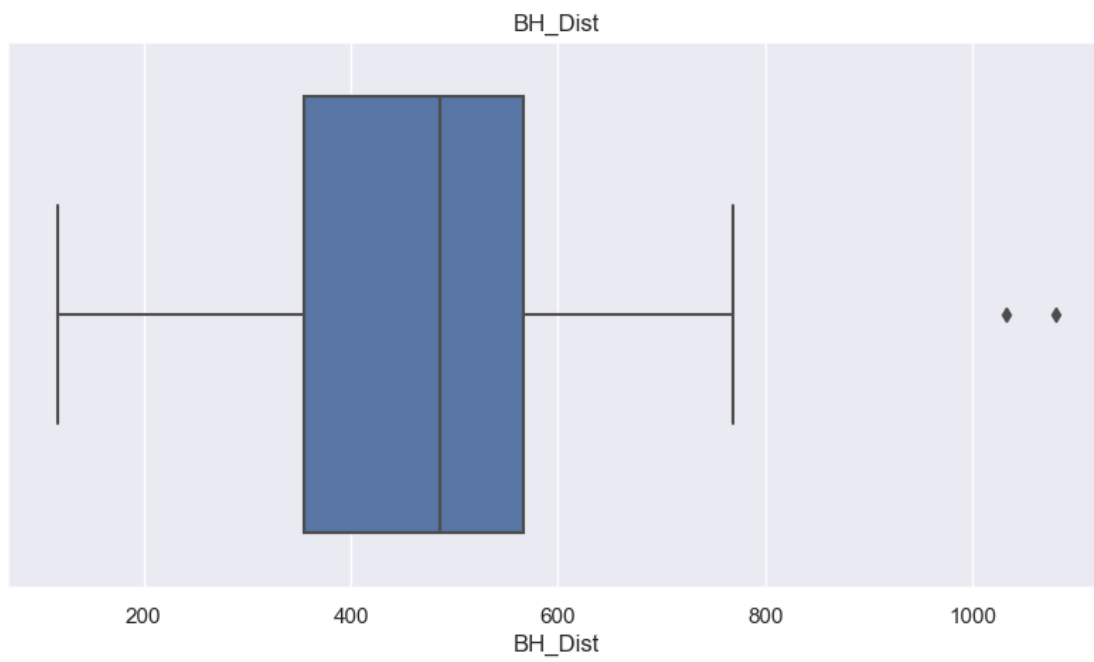
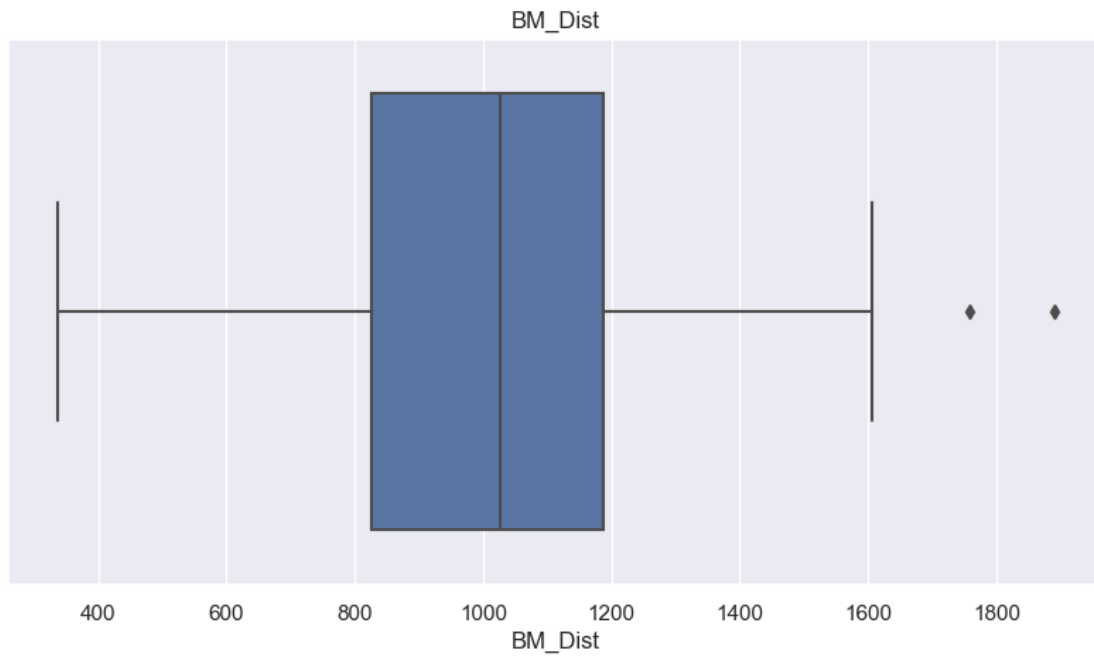


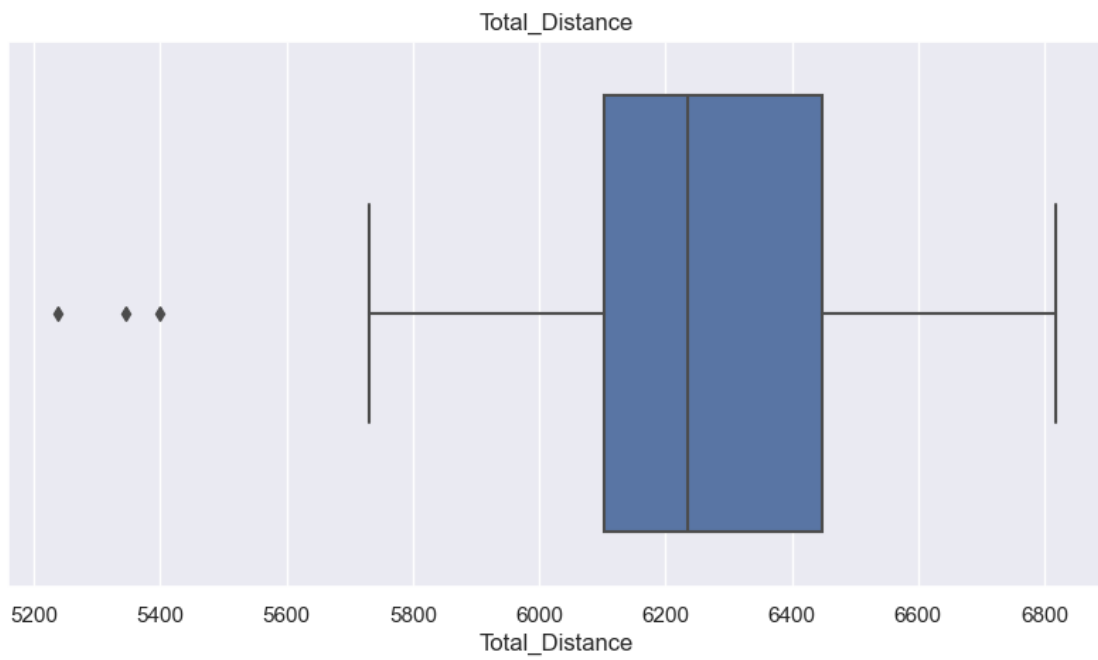
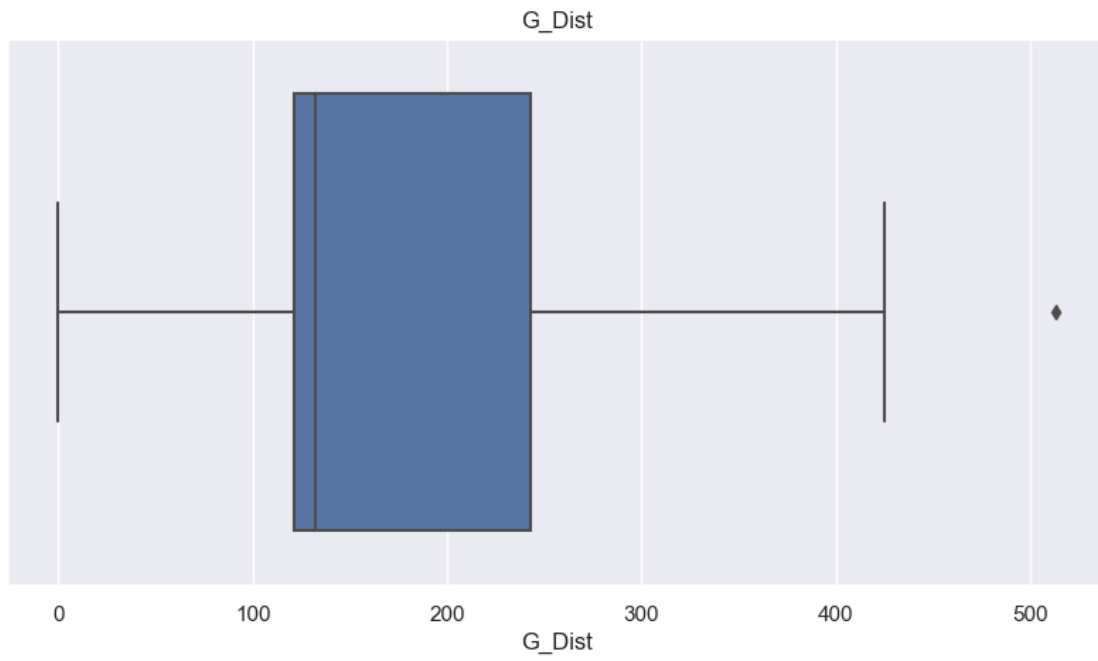
```
[5]: # Boxplot for all features to identify outliers
for column in df[['Target', 'Myelin', 'D_Dist',
    ↳ 'T_Dist', 'A_Dist', 'BL_Dist', 'BM_Dist', 'BH_Dist', 'G_Dist', 'Total_Distance', 'Total_Rate', 'Web
    ↳ columns:
    if df[column].dtype in ['float64', 'int64']: # only plot for numeric
    ↳ columns
        plt.figure(figsize=(10, 5))
        sns.boxplot(x=column, data=df)
        plt.title(column)
        plt.show()
```

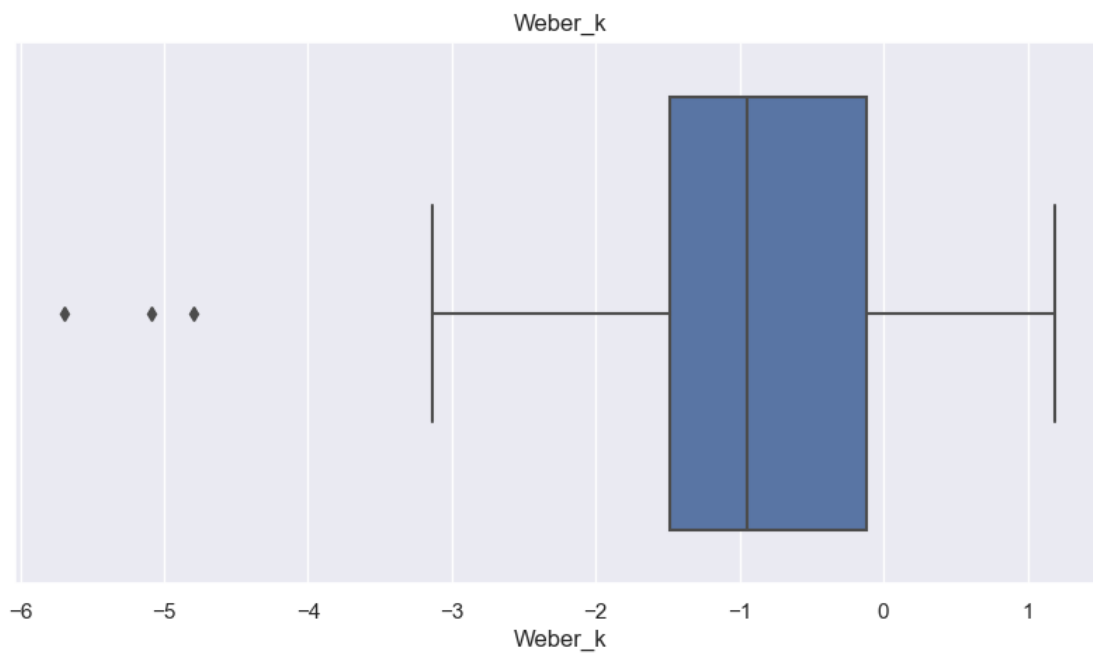
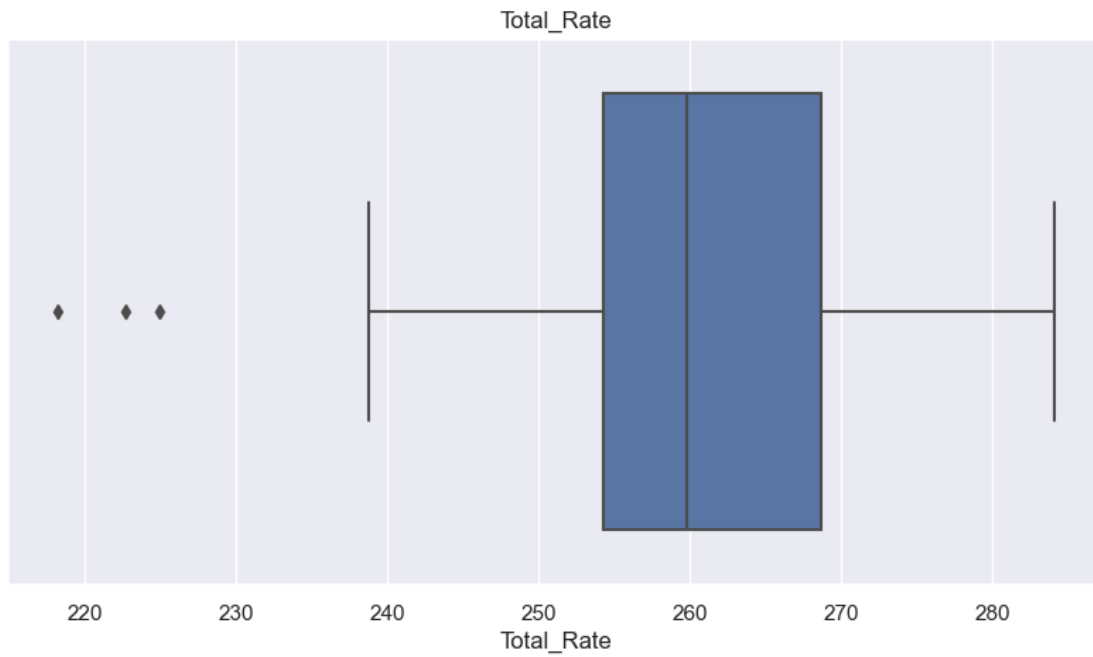




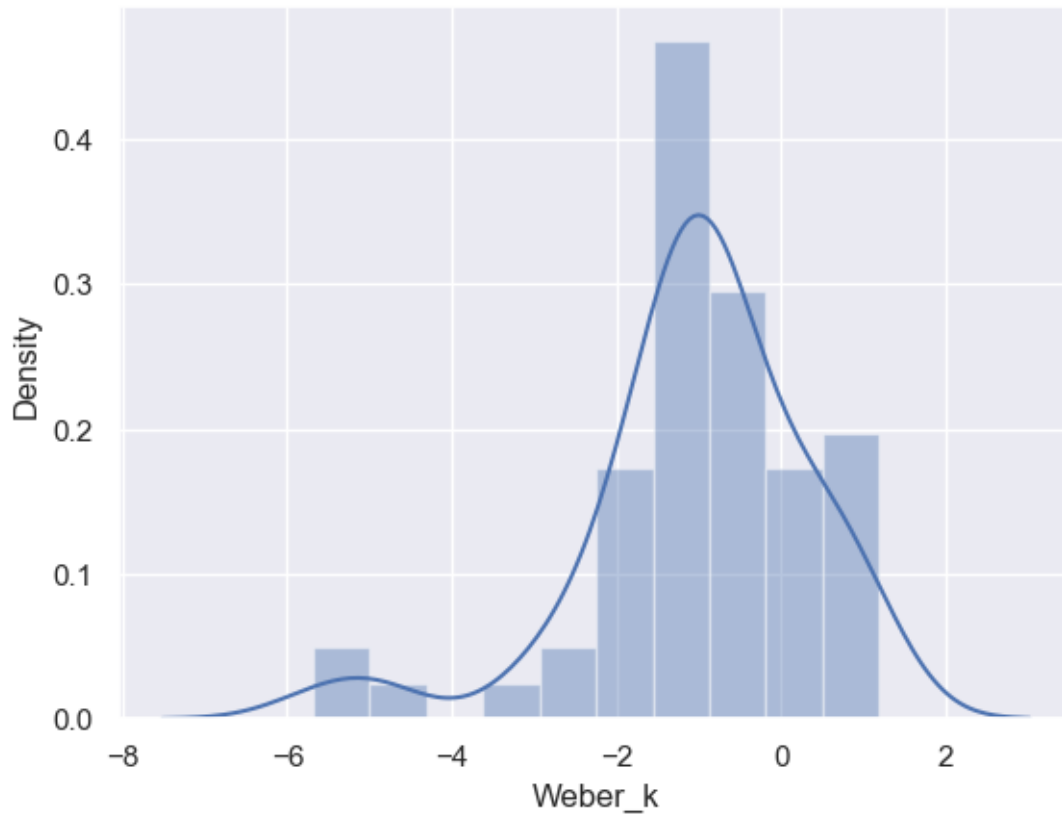








```
[6]: # Distribution of Weber constant  
_ = sns.distplot(df.Weber_k)
```



```
[7]: # Select features and define a subset
selected_columns = ['Target', 'Myelin', 'D_Dist', 'T_Dist', 'A_Dist',
                    ↪ 'BL_Dist', 'BM_Dist', 'BH_Dist', 'G_Dist', 'Total_Distance', 'Total_Rate',
                    ↪ 'Weber_k']

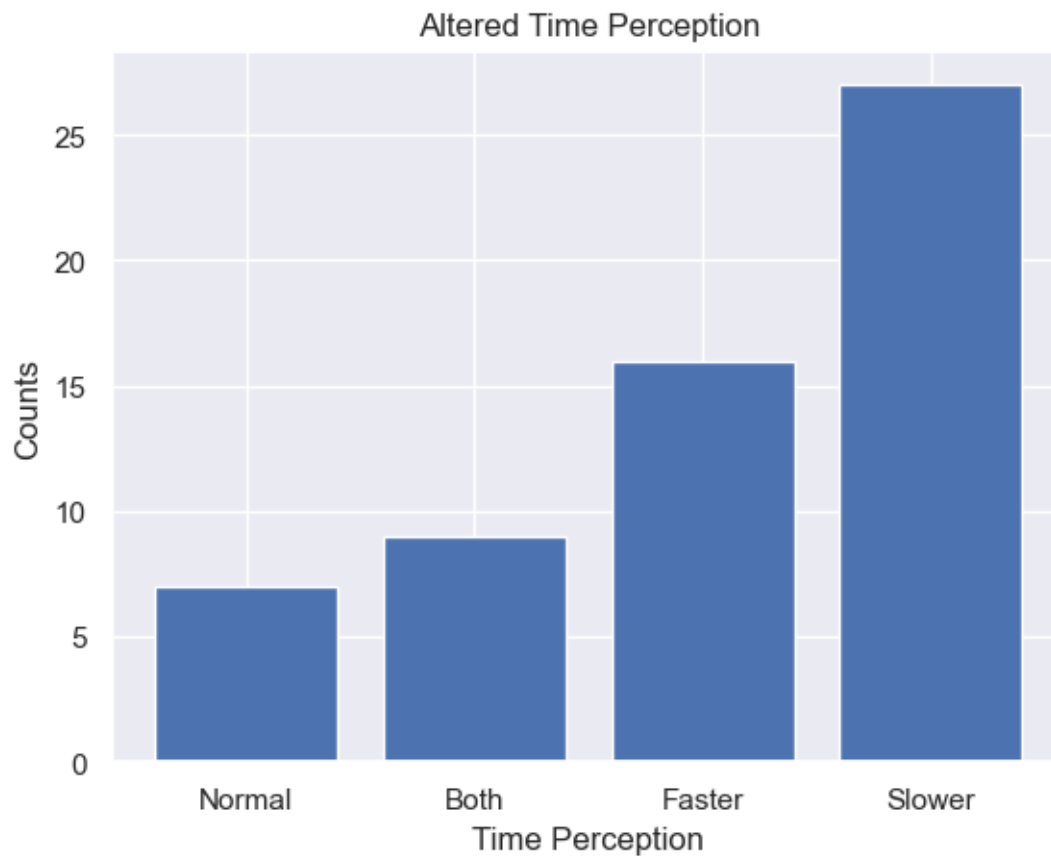
subset_df = df[selected_columns]
```

```
[8]: # Label Mapping
label_mapping = {
    3: 'Slower',
    2: 'Faster',
    1: 'Both',
    0: 'Normal'
}

# Get value counts
counts = df['Target'].value_counts().sort_index()

# Use the labels from label_mapping for plotting
labels = [label_mapping[key] for key in counts.index]
```

```
# Plot the value counts
plt.bar(labels, counts)
plt.title('Altered Time Perception')
plt.xlabel('Time Perception')
plt.ylabel('Counts')
plt.show()
```



```
[9]: # Examine Target value counts
subset_df['Target'].value_counts()
```

```
[9]: 3    27
     2    16
     1     9
     0     7
     Name: Target, dtype: int64
```

8 Subset Models

A Decision Tree classifier is a versatile machine learning algorithm that offers high interpretability by visualizing decision-making pathways. It can handle both numerical and categorical data, capture non-linear relationships, and requires no feature scaling.

```
[10]: # Split the data into training and test sets
X = subset_df.drop('Target', axis=1)
y = subset_df['Target']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
    random_state=42)

# Initialize a Decision Tree classifier.
# Decision Trees are a non-parametric supervised learning method used for both
# classification and regression. It works by partitioning the source set into
# subsets based on the values of input attributes.
clf = DecisionTreeClassifier()

# Train the Decision Tree classifier on the training data.
# The fit method will construct a tree from the training data, trying to split
# on features and make decisions in a way that accurately predicts the target
# variable.
clf.fit(X_train, y_train)

# Once the model is trained, use it to predict the target variable for the test
    data.
# The predict method traverses the trained decision tree to produce a
    prediction
# for each test sample.
y_pred = clf.predict(X_test)

# Defining mapping of numeric labels to their corresponding word labels
numeric_labels = [0, 1, 2, 3]
word_labels = ["Average", "Both", "Faster", "Slower"]

# Check accuracy
accuracy = accuracy_score(y_test, y_pred)*100
print()
print(f"Accuracy: {accuracy:.2f}%")
print()
print(classification_report(y_test, y_pred, target_names=word_labels))

# Feature names
feature_names = [
    'Myelin', 'D_Dist', 'T_Dist', 'A_Dist', 'BL_Dist',
    'BM_Dist', 'BH_Dist', 'G_Dist', 'Total_Distance', 'Total_Rate', 'Weber_k'
]
```

```

# Extracting feature importances from the Decision Tree model
feature_importances = clf.feature_importances_

# Pairing feature names with their importances and sorting them
sorted_importances = sorted(zip(feature_names, feature_importances), key=lambda_
    ↪x: x[1], reverse=True)

# Display
print("\nFeature Importances:")
print()
for feature, importance in sorted_importances:
    print(f"{feature}: {importance:.4f}")

# Define numeric labels and corresponding word labels
numeric_labels = [0, 1, 2, 3]
word_labels = ["Average", "Both", "Faster", "Slower"]

# Create a confusion matrix
cm = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
                    ↪color=color)
            else:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
                    ↪color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

# Use the existing plot function

```

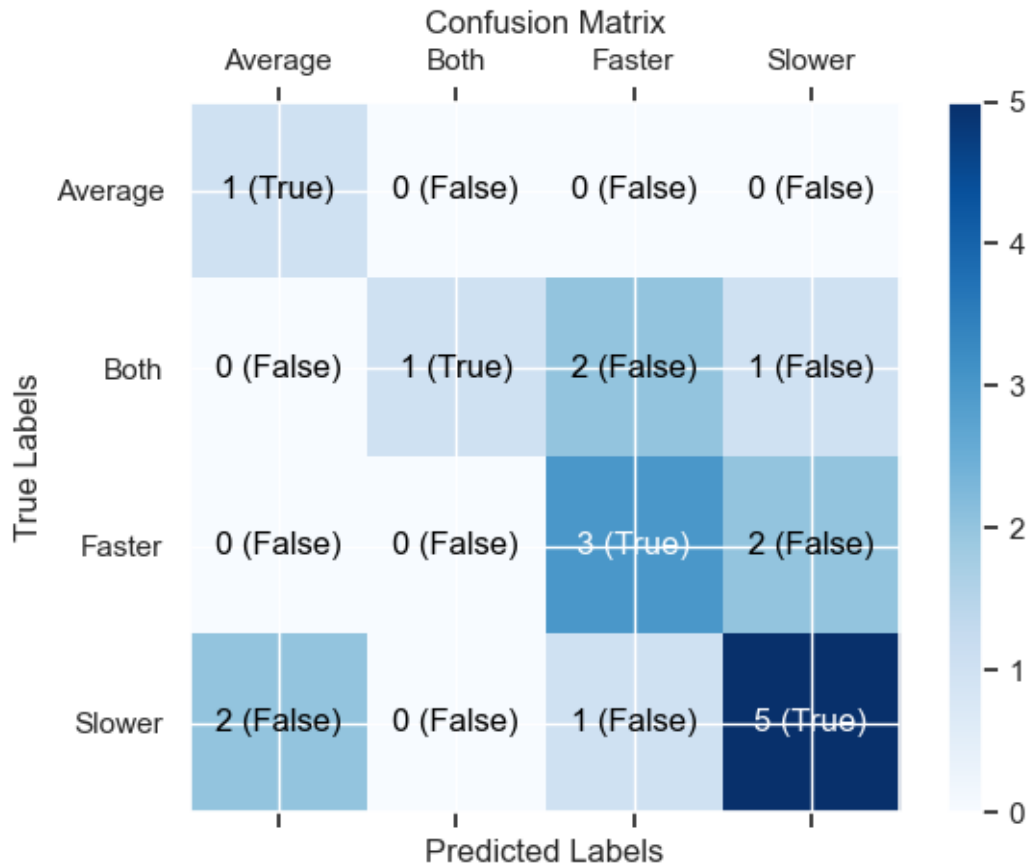
```
plot_confusion_matrix(cm, word_labels)
```

Accuracy: 55.56%

	precision	recall	f1-score	support
Average	0.33	1.00	0.50	1
Both	1.00	0.25	0.40	4
Faster	0.50	0.60	0.55	5
Slower	0.62	0.62	0.62	8
accuracy			0.56	18
macro avg	0.61	0.62	0.52	18
weighted avg	0.66	0.56	0.55	18

Feature Importances:

BM_Dist: 0.3252
BL_Dist: 0.2101
Total_Distance: 0.1249
D_Dist: 0.1201
BH_Dist: 0.0902
G_Dist: 0.0576
Total_Rate: 0.0480
A_Dist: 0.0238
Myelin: 0.0000
T_Dist: 0.0000
Weber_k: 0.0000



```
[11]: # Feature names for your dataset
feature_names = ['Myelin', 'D_Dist', 'T_Dist', 'A_Dist', 'BL_Dist', 'BM_Dist', 'BH_Dist', 'G_Dist', 'Total_Distance', 'Total_Rate', 'Weber_k']

# Export the Decision Tree to a dot format
dot_data = export_graphviz(clf, out_file=None,
                           feature_names=feature_names,
                           class_names=word_labels,
                           filled=True, rounded=True,
                           special_characters=True)

# Use graphviz to create the graph object
graph = graphviz.Source(dot_data)

# View the graph
graph.view(filename="C:/Users/newmy/Desktop/Temporal_Metrics/DT3", cleanup=True)

# Render and save the graph to a file (this is technically redundant after the
# view command, but ensures the file is saved)
```

```
graph.render(filename="C:/Users/newmy/Desktop/Temporal_Metrics/DT3",  
             ↪format='pdf', cleanup=True)
```

```
[11]: 'C:\\Users\\newmy\\Desktop\\Temporal_Metrics\\DT3.pdf'
```

The following code is setting up and training a Random Forest classifier. A Random Forest is an ensemble learning method that creates a ‘forest’ of decision trees during training and outputs the class that is the mode of the classes (classification) or mean/average prediction (regression) of the individual trees for a given input.

```
[12]: # Extracting features and target variable from the dataset  
X = subset_df.drop('Target', axis=1) # Features (excluding the target variable)  
y = subset_df['Target'] # Target variable  
  
# Splitting the data into training and testing sets, with 30% of the data being  
↪used for testing  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,  
            ↪random_state=42)  
  
# Initializing a Random Forest classifier with 100 trees  
clf = RandomForestClassifier(n_estimators=100, random_state=42)  
  
# Training the Random Forest classifier on the training data  
clf.fit(X_train, y_train)  
  
# Predicting the target variable for the testing set  
y_pred = clf.predict(X_test)  
  
# Defining mapping of numeric labels to their corresponding word labels  
numeric_labels = [0, 1, 2, 3]  
word_labels = ["Average", "Both", "Faster", "Slower"]  
  
accuracy_percentage = accuracy_score(y_test, y_pred) * 100  
print(f"Accuracy: {accuracy_percentage:.2f}%")  
print()  
print(classification_report(y_test, y_pred, target_names=word_labels))  
  
# Extract the feature importances  
feature_importances = clf.feature_importances_  
  
# Combine feature names and their importance scores  
features_df = pd.DataFrame({  
    'Feature': X_train.columns,  
    'Importance': feature_importances  
})  
  
# Sort by importance
```



```

features_df = features_df.sort_values(by='Importance', ascending=False)

print(features_df.to_string(index=False))

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
                    color=color)
            else:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
                    color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

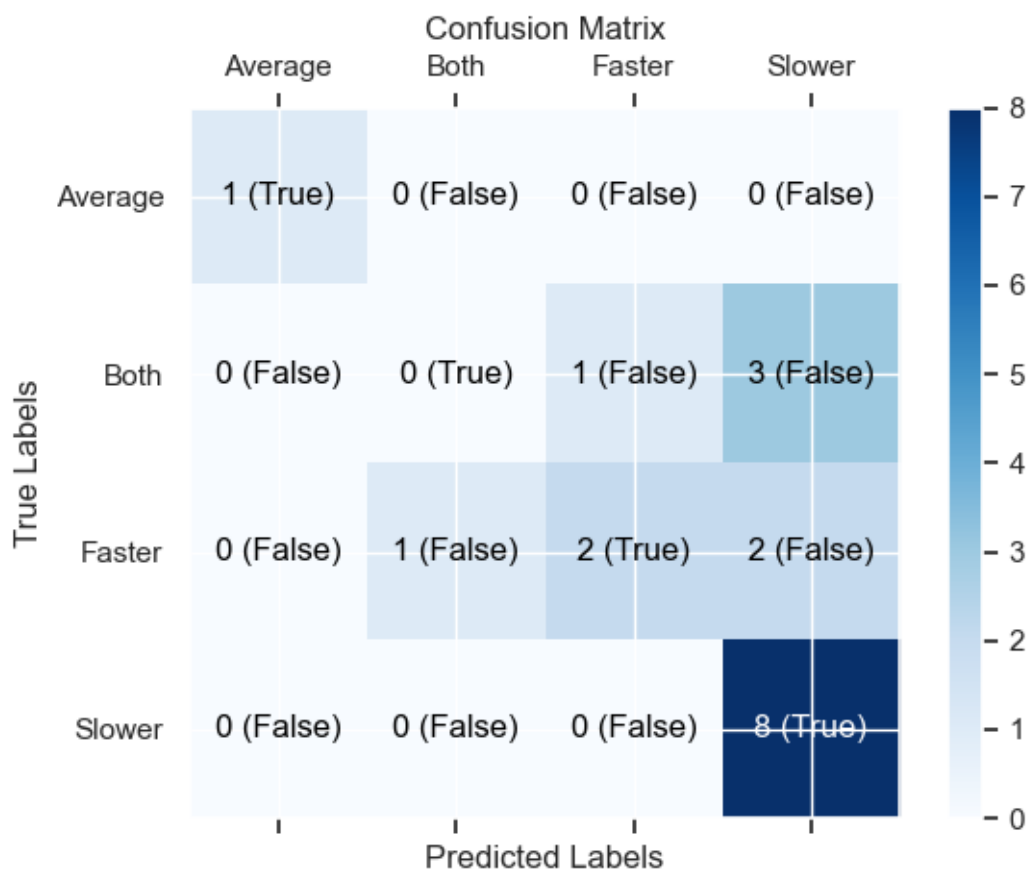
# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)

```

Accuracy: 61.11%

	precision	recall	f1-score	support
Average	1.00	1.00	1.00	1
Both	0.00	0.00	0.00	4
Faster	0.67	0.40	0.50	5
Slower	0.62	1.00	0.76	8
accuracy			0.61	18
macro avg	0.57	0.60	0.57	18
weighted avg	0.51	0.61	0.53	18

Feature	Importance
BM_Dist	0.148675
BL_Dist	0.143355
Weber_k	0.101638
Total_Rate	0.092264
BH_Dist	0.091719
Total_Distance	0.087851
D_Dist	0.076627
Myelin	0.070906
G_Dist	0.063853
A_Dist	0.061914
T_Dist	0.061196



The following code is setting up and training a Random Forest classifier. However, unlike the previous code, it also includes a preprocessing step for standardizing features and uses a grid search to optimize hyperparameters for the Random Forest model.

```
[13]: # Extracting features and target variable from the dataset
X = subset_df.drop('Target', axis=1) # Features (excluding the target variable)
y = subset_df['Target'] # Target variable

# Splitting the data into training and testing sets. We use stratify to ensure
# the training and test datasets have similar proportion of target classes.
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
    ↪random_state=42, stratify=y)

# Standardizing the features. This step is optional for Random Forests since
    ↪they are
# scale invariant, but is included for demonstration purposes.
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Initializing a base Random Forest classifier model
rf = RandomForestClassifier(random_state=42)

# Defining a grid of hyperparameters to optimize the Random Forest model
param_grid = {
    'n_estimators': [50, 100, 200],
    'max_depth': [None, 10, 20, 30],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4],
    'bootstrap': [True, False]
}

# Using GridSearchCV to search for the best hyperparameters over the specified
    ↪grid.
# The search will be based on 3-fold cross-validation and will use all CPU
    ↪cores (n_jobs=-1).
grid_search = GridSearchCV(estimator=rf, param_grid=param_grid,
    cv=3, n_jobs=-1, verbose=2, scoring='accuracy')

# Training the model using GridSearchCV to find the best hyperparameters
grid_search.fit(X_train, y_train)

# Extracting the best Random Forest model after grid search
best_rf = grid_search.best_estimator_

# Predicting target variable for the test set using the best model
y_pred = best_rf.predict(X_test)

# Evaluating the model
accuracy_percentage = accuracy_score(y_test, y_pred) * 100
```

```

print()
print(f"Accuracy: {accuracy_percentage:.2f}%")
print()
print(classification_report(y_test, y_pred, target_names=word_labels))

# Extract the feature importances
feature_importances = best_rf.feature_importances_

# Combine feature names and their importance scores
features_df = pd.DataFrame({
    'Feature': X.columns,
    'Importance': feature_importances
})

# Sort by importance
features_df = features_df.sort_values(by='Importance', ascending=False)

print(features_df.to_string(index=False))

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
                    ↪color=color)
            else:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
                    ↪color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

```

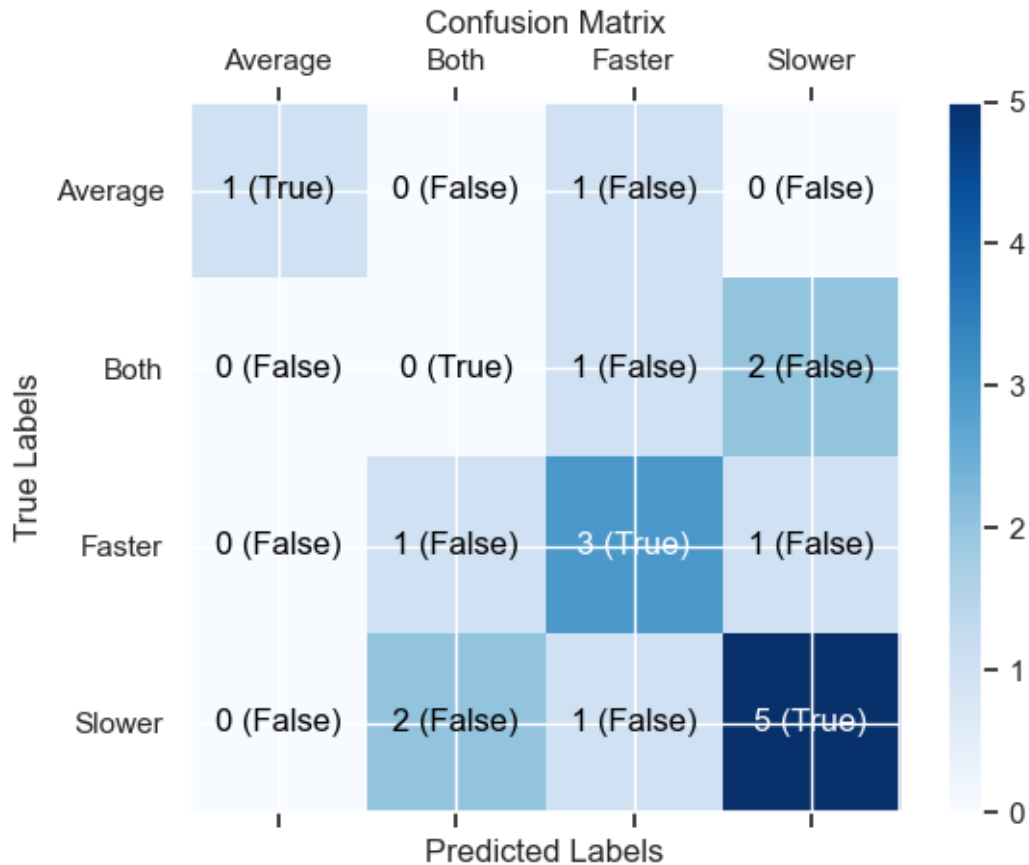
```
# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)
```

Fitting 3 folds for each of 216 candidates, totalling 648 fits

Accuracy: 50.00%

	precision	recall	f1-score	support
Average	1.00	0.50	0.67	2
Both	0.00	0.00	0.00	3
Faster	0.50	0.60	0.55	5
Slower	0.62	0.62	0.62	8
accuracy			0.50	18
macro avg	0.53	0.43	0.46	18
weighted avg	0.53	0.50	0.50	18

Feature	Importance
BH_Dist	0.162598
BM_Dist	0.158564
G_Dist	0.125259
Myelin	0.108026
Total_Distance	0.091371
BL_Dist	0.086880
Total_Rate	0.072096
Weber_k	0.067589
D_Dist	0.049068
A_Dist	0.048514
T_Dist	0.030035



This code is preparing a dataset, splitting it into training and test subsets, pre-processing the data, initializing a Gradient Boosting Classifier, training it on the training data, and finally using the trained model to predict the target variable for the test data.

```
[14]: # Drop the 'Target' column to get the feature matrix 'X'
X = subset_df.drop('Target', axis=1)

# Isolate the 'Target' column to get the target vector 'y'
y = subset_df['Target']

# Splitting the dataset into training and test sets.
# Using stratify ensures that the distribution of the target variable is
↳ consistent between training and test sets.
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
↳ random_state=42, stratify=y)

# Initialize a standard scaler. This will mean-center and scale the data to
↳ unit variance.
```

```

# This step isn't necessary for Gradient Boosting since it's based on decision
# trees, but
# sometimes it's used to maintain a consistent pre-processing pipeline or to
# help with convergence.
scaler = StandardScaler()

# Fit the scaler on the training data and transform it.
X_train = scaler.fit_transform(X_train)

# Transform the test data using the same scaler (no fitting here to prevent
# data leakage).
X_test = scaler.transform(X_test)

# Initialize a Gradient Boosting Classifier.
# n_estimators represents the number of boosting stages to be run.
# learning_rate shrinks the contribution of each tree.
# max_depth is the maximum depth of the individual trees.
gb = GradientBoostingClassifier(n_estimators=100, learning_rate=0.1,
# max_depth=3, random_state=42)

# Train the Gradient Boosting Classifier on the training data.
gb.fit(X_train, y_train)

# Use the trained Gradient Boosting model to make predictions on the test data.
y_pred = gb.predict(X_test)

# Evaluating the model
accuracy_percentage = accuracy_score(y_test, y_pred) * 100
print(f"Accuracy: {accuracy_percentage:.2f}%")
print()

print(classification_report(y_test, y_pred, target_names=word_labels))

# Extract the feature importances
feature_importances = clf.feature_importances_

# Combine feature names and their importance scores
features_df = pd.DataFrame({
    'Feature': X.columns,
    'Importance': feature_importances
})

# Sort by importance
features_df = features_df.sort_values(by='Importance', ascending=False)

print(features_df.to_string(index=False))

```

```

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
                    ↪color=color)
            else:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
                    ↪color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)

```

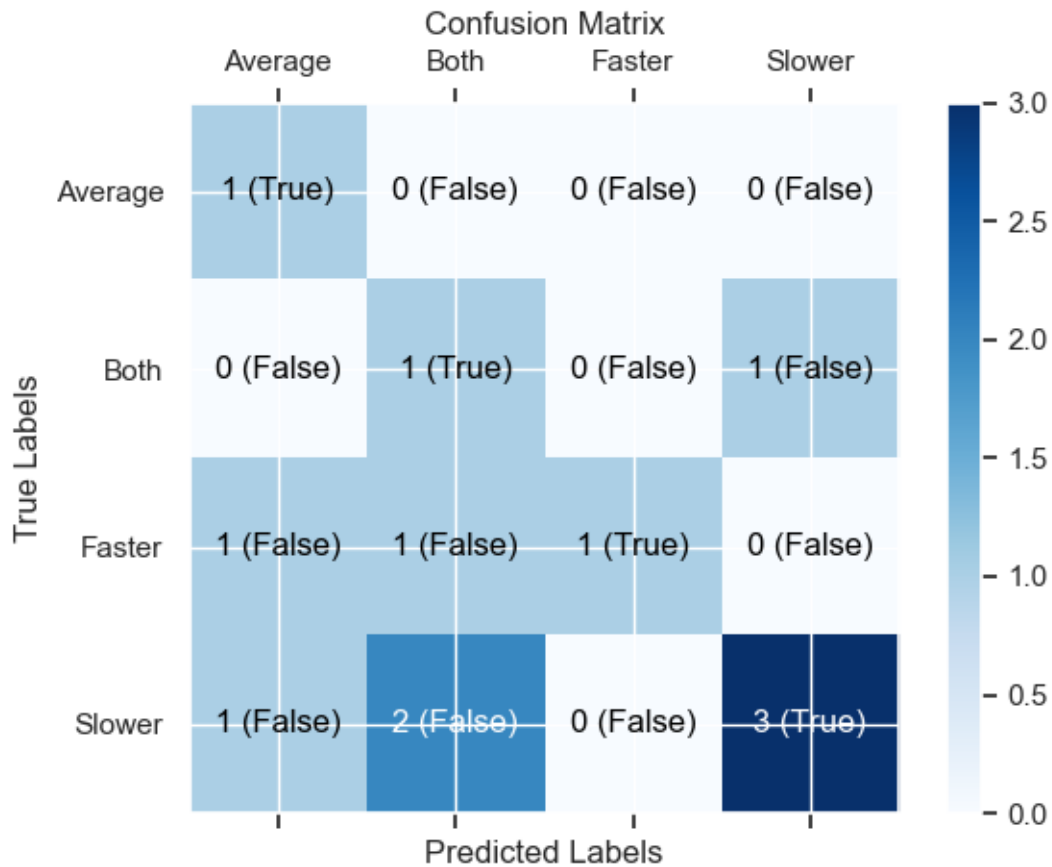
Accuracy: 50.00%

	precision	recall	f1-score	support
Average	0.33	1.00	0.50	1
Both	0.25	0.50	0.33	2
Faster	1.00	0.33	0.50	3
Slower	0.75	0.50	0.60	6
accuracy			0.50	12
macro avg	0.58	0.58	0.48	12
weighted avg	0.69	0.50	0.52	12
Feature	Importance			
BM_Dist	0.148675			
BL_Dist	0.143355			


```

Weber_k      0.101638
Total_Rate   0.092264
BH_Dist      0.091719
Total_Distance 0.087851
D_Dist       0.076627
Myelin       0.070906
G_Dist       0.063853
A_Dist       0.061914
T_Dist       0.061196

```



This code demonstrates how to handle imbalanced datasets using SMOTE to generate synthetic instances of the minority class, and then trains a Random Forest Classifier on the balanced training set to make predictions on the test data.

```

[15]: # Extract features by dropping the 'Target' column, resulting in a features_
      ↪matrix 'X'
X = subset_df.drop("Target", axis=1)

# Extract the 'Target' column to form the target vector 'y'
y = subset_df["Target"]

```

```

# Split the dataset into a training subset and a test subset
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

# SMOTE (Synthetic Minority Over-sampling Technique) is an over-sampling method
↳ that creates synthetic examples
# in the feature space. It's used to handle imbalanced datasets by increasing
↳ the number of instances in the minority class.
sm = SMOTE(k_neighbors=3) # Using 3 nearest neighbors

# Apply SMOTE to the training data. This results in a balanced (or more
↳ balanced) training dataset.
X_train_resampled, y_train_resampled = sm.fit_resample(X_train, y_train)

# Initialize a Random Forest Classifier. Random Forest is an ensemble learning
↳ method
# that constructs multiple decision trees during training and outputs the
↳ majority class
# (for classification problems) of the individual trees for predictions.
clf = RandomForestClassifier(n_estimators=100) # Using 100 trees in the forest

# Train the Random Forest Classifier on the resampled (balanced) training data
clf.fit(X_train_resampled, y_train_resampled)

# Use the trained Random Forest model to predict the target for the test data
y_pred = clf.predict(X_test)

# Check accuracy
accuracy_percentage = accuracy_score(y_test, y_pred) * 100
print(f"Accuracy: {accuracy_percentage:.2f}%")
print()
print(classification_report(y_test, y_pred, target_names=word_labels))
print()
# Extract the feature importances
feature_importances = clf.feature_importances_

# Combine feature names and their importance scores
features_df = pd.DataFrame({
    'Feature': X.columns,
    'Importance': feature_importances
})

# Sort by importance
features_df = features_df.sort_values(by='Importance', ascending=False)

print(features_df.to_string(index=False))
print()

```

```

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
                    ↪color=color)
            else:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
                    ↪color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)

```

Accuracy: 55.56%

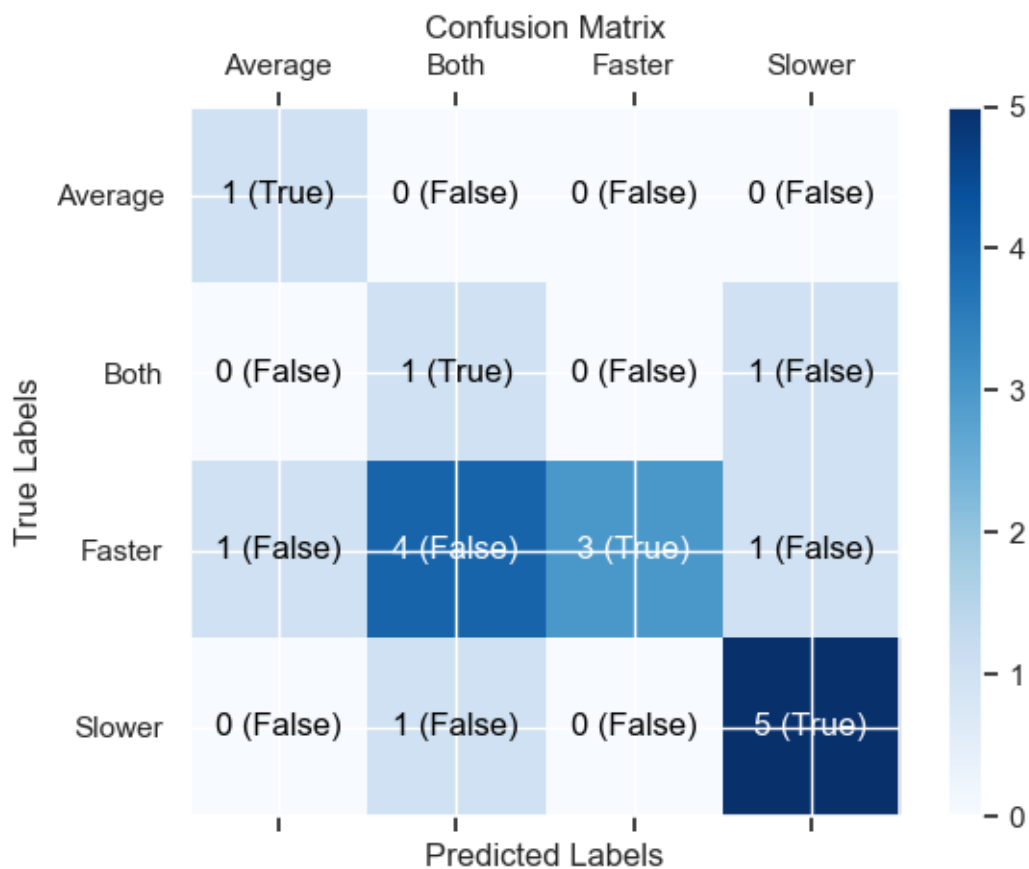
	precision	recall	f1-score	support
Average	0.50	1.00	0.67	1
Both	0.17	0.50	0.25	2
Faster	1.00	0.33	0.50	9
Slower	0.71	0.83	0.77	6
accuracy			0.56	18
macro avg	0.60	0.67	0.55	18
weighted avg	0.78	0.56	0.57	18

Feature Importance

```

BM_Dist      0.137574
BL_Dist      0.102968
Myelin       0.100807
Weber_k      0.090209
BH_Dist      0.089271
D_Dist       0.086229
T_Dist       0.085837
Total_Rate   0.083863
G_Dist       0.077844
A_Dist       0.075564
Total_Distance 0.069834

```



In essence, this code demonstrates the process of training a Random Forest model, extracting feature importances to identify which features are the most informative, and then re-training the model using only those top features to make predictions on the test data.

```

[16]: # Initialize a Random Forest Classifier. Random Forest is an ensemble learning
      ↪method

```

```

# that constructs multiple decision trees during training and outputs the
↳majority class
# (for classification problems) of the individual trees for predictions.
clf = RandomForestClassifier(n_estimators=100) # Using 100 trees in the forest

# Train the Random Forest Classifier on the training data
clf.fit(X_train, y_train)

# Obtain feature importances from the trained Random Forest model. This gives
↳insight
# into which features the model found to be the most informative for making
↳predictions.
feature_importances = clf.feature_importances_

# Convert the feature importances to a DataFrame for easier visualization and
↳sorting
features_df = pd.DataFrame({
    'Feature': X.columns,          # Feature names
    'Importance': feature_importances # Their corresponding importance scores
})

# Display sorted feature importances
print("Feature Importances:")
print()
sorted_features = features_df.sort_values(by="Importance", ascending=False)
print(sorted_features.to_string(index=False)) # Display without the default
↳index for cleaner output

# Based on the sorted importances, select the top features.
# Here, we're assuming we want the top 10 most important features.
top_features = sorted_features['Feature'].head(10).tolist()

# Subset the training and test data to include only these top features
X_train_selected = X_train[top_features]
X_test_selected = X_test[top_features]

# Train the Random Forest Classifier using only the top features. This might
↳lead to a more focused and
# possibly better performing model if some features were not informative or
↳were introducing noise.
clf.fit(X_train_selected, y_train)

# Use the re-trained Random Forest model to predict the target for the test data
y_pred = clf.predict(X_test_selected)

# Check accuracy

```

```

accuracy_percentage = accuracy_score(y_test, y_pred) * 100
print()
print(f"Accuracy: {accuracy_percentage:.2f}%")
print()
print(classification_report(y_test, y_pred, target_names=word_labels))
print()
# Define numeric labels and corresponding word labels
numeric_labels = [0, 1, 2, 3]
word_labels = ["Average", "Both", "Faster", "Slower"]

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
                    ↪color=color)
            else:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
                    ↪color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)

```

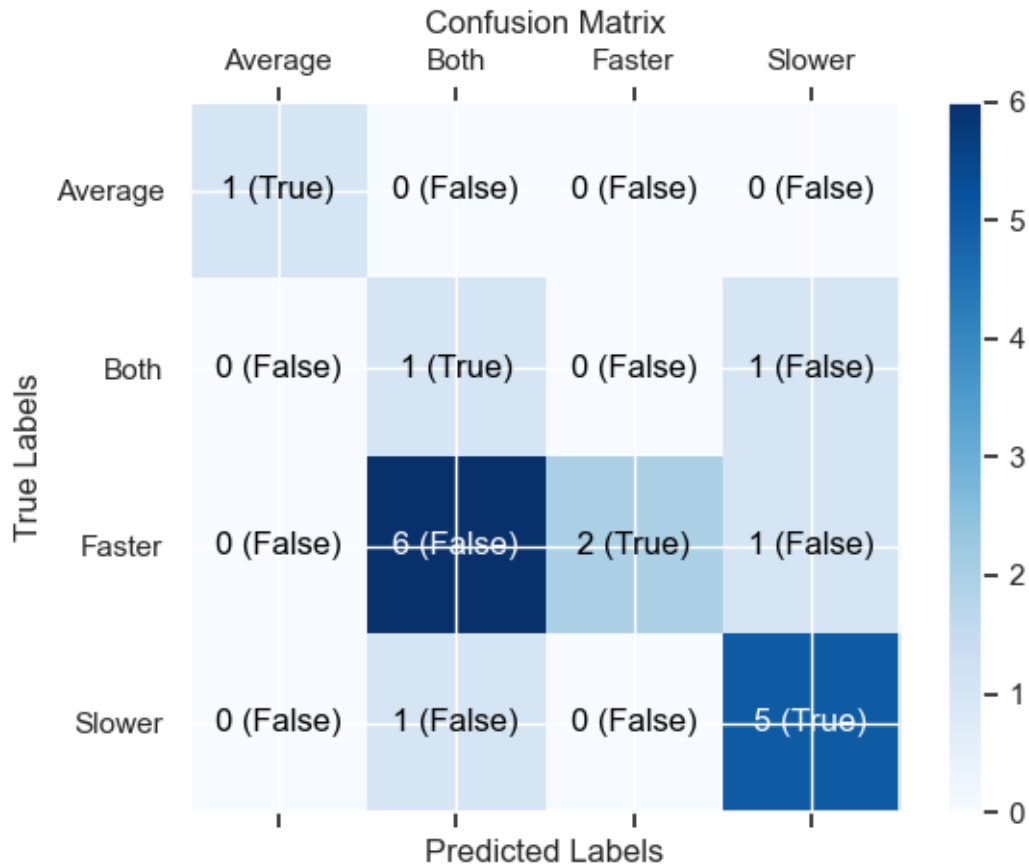
Feature Importances:

Feature	Importance
BM_Dist	0.165945
BL_Dist	0.114394
BH_Dist	0.105832

Total_Distance	0.091792
Total_Rate	0.086174
Weber_k	0.082808
Myelin	0.077146
G_Dist	0.076108
T_Dist	0.075753
D_Dist	0.064607
A_Dist	0.059442

Accuracy: 50.00%

	precision	recall	f1-score	support
Average	1.00	1.00	1.00	1
Both	0.12	0.50	0.20	2
Faster	1.00	0.22	0.36	9
Slower	0.71	0.83	0.77	6
accuracy			0.50	18
macro avg	0.71	0.64	0.58	18
weighted avg	0.81	0.50	0.52	18



This code constructs an ensemble classifier, called a “Voting Classifier”, using three individual models: a Random Forest, an SVM (Support Vector Machine), and a Logistic Regression. The ensemble classifier takes the predictions of each individual model and aggregates them to produce a final prediction. After training, the code computes the ensemble classifier’s accuracy on test data.

```
[17]: # Initialize individual models.

# RandomForest is an ensemble learning method that constructs multiple decision
# trees
# during training and outputs the majority class (for classification problems)
# of
# the individual trees for predictions.
clf1 = RandomForestClassifier(n_estimators=100)

# SVM is a supervised machine learning algorithm which can be used for
# classification
# or regression problems. It uses a technique called the kernel trick to
# transform
```



```

# your data and then based on these transformations it finds an optimal
↳ boundary
# between the possible outputs.
clf2 = SVC(probability=True) # probability=True allows to obtain probabilities
↳ with predict_proba

# Logistic Regression is a statistical model that in its basic form uses a
↳ logistic function
# to model a binary dependent variable.
clf3 = LogisticRegression()

# Create a voting classifier that combines the predictions from the three
↳ individual models.
# The 'soft' voting means predictions are based on argmax of the sums of the
↳ predicted
# probabilities, which recommends weighted average.
ecf = VotingClassifier(estimators=[
    ('rf', clf1), ('svc', clf2), ('lr', clf3)], voting='soft')

# Train the ensemble model on training data.
ecf.fit(X_train, y_train)

# Use the trained ensemble model to predict the target variable on the test
↳ data.
y_pred = ecf.predict(X_test)

# Evaluate the ensemble model's performance.

# Calculate the accuracy of the ensemble model on the test data.
accuracy_ensemble = accuracy_score(y_test, y_pred)*100
print()
print(f"Ensemble Accuracy: {accuracy_ensemble:.2f}%")
print()

# Print a classification report showing various metrics (precision, recall,
↳ f1-score, etc.)
classification_rep_ensemble = classification_report(y_test, y_pred,
↳ target_names=word_labels)
print(classification_rep_ensemble)

# Extracting and combining feature importances from the models.

# List of feature names for reference.
feature_names = [
    'Myelin', 'D_Dist', 'T_Dist', 'A_Dist', 'BL_Dist',
    'BM_Dist', 'BH_Dist', 'G_Dist', 'Total_Distance', 'Total_Rate', 'Weber_k'

```

```

]

# Access each individual model inside the VotingClassifier after training.
fitted_rf = eclf.named_estimators_['rf']
fitted_svc = eclf.named_estimators_['svc']
fitted_lr = eclf.named_estimators_['lr']

# Extract feature importances from RandomForest.
feature_importances_rf = fitted_rf.feature_importances_

# Extract feature importances from SVM.
# Note: Importances from SVM are relevant only when the SVM uses a linear
↳ kernel.
if isinstance(fitted_svc.kernel, str) and fitted_svc.kernel == 'linear':
    feature_importances_svc = abs(fitted_svc.coef_[0])
else:
    feature_importances_svc = np.ones(len(feature_names)) # Assign uniform
↳ importance for non-linear SVM.

# Extract feature importances from Logistic Regression.
feature_importances_lr = abs(fitted_lr.coef_[0])

# Helper function to normalize the feature importances.
def normalize(importance):
    return importance / sum(importance)

# Normalize the feature importances for each model.
normalized_rf = normalize(feature_importances_rf)
normalized_svc = normalize(feature_importances_svc)
normalized_lr = normalize(feature_importances_lr)

# Combine (sum) normalized importances from all models.
combined_importance = normalized_rf + normalized_svc + normalized_lr

# Pair the feature names with their combined importances and sort them in
↳ descending order.
sorted_importances = sorted(zip(feature_names, combined_importance), key=lambda
↳ x: x[1], reverse=True)
# Display
print("Feature Importances:")
print()
for feature, importance in sorted_importances:
    print(f"{feature}: {importance:.4f}")

# Define numeric labels and corresponding word labels
numeric_labels = [0, 1, 2, 3]
word_labels = ["Average", "Both", "Faster", "Slower"]

```

```

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
                    ↪color=color)
            else:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
                    ↪color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)

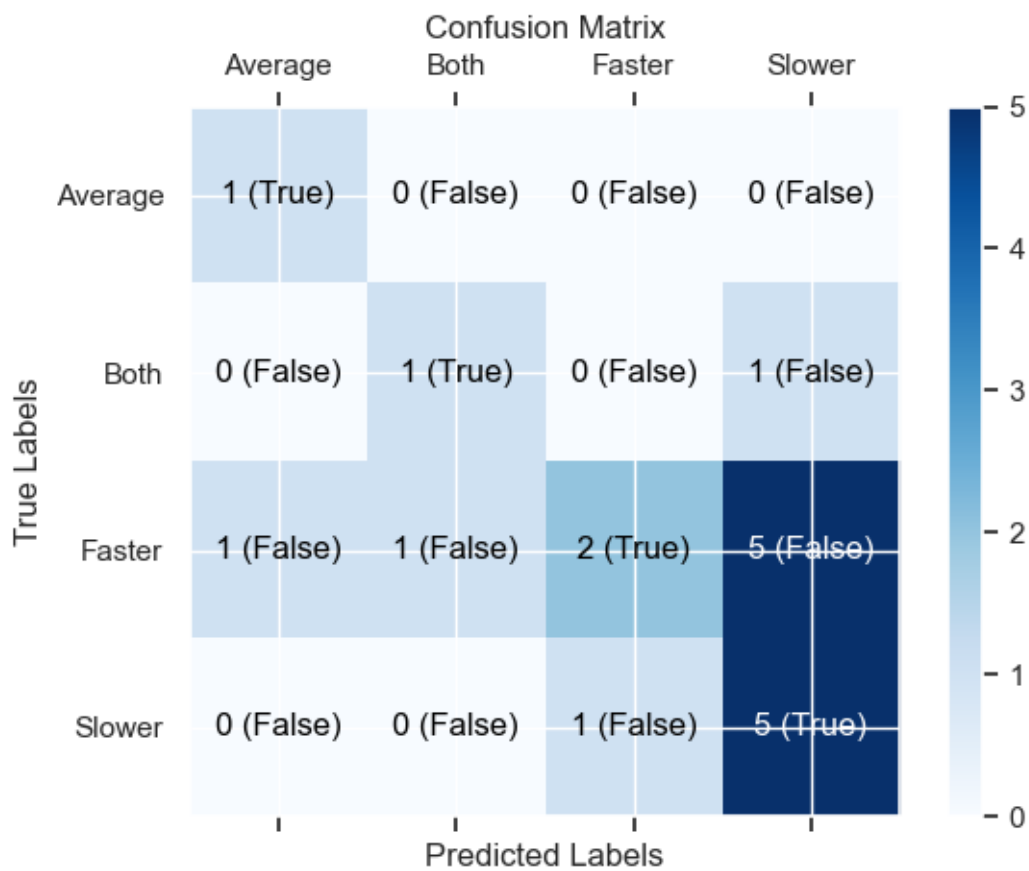
```

Ensemble Accuracy: 50.00%

	precision	recall	f1-score	support
Average	0.50	1.00	0.67	1
Both	0.50	0.50	0.50	2
Faster	0.67	0.22	0.33	9
Slower	0.45	0.83	0.59	6
accuracy			0.50	18
macro avg	0.53	0.64	0.52	18
weighted avg	0.57	0.50	0.46	18

Feature Importances:

BH_Dist: 0.4580
 T_Dist: 0.3834
 BM_Dist: 0.3651
 G_Dist: 0.3501
 D_Dist: 0.2498
 A_Dist: 0.2446
 BL_Dist: 0.2396
 Weber_k: 0.1986
 Total_Distance: 0.1765
 Total_Rate: 0.1763
 Myelin: 0.1581



9 Full Dataset Models

```
[18]: feature_names = df.columns.tolist()
      print(feature_names)
```

```
['Condition', 'Target', 'Myelin', 'Speed', 'Speed_M', 'D_WL', 'D_Cycle',
```

```
'D_Dist', 'D_Time', 'D_Freq', 'T_WL', 'T_Cycle', 'T_Dist', 'T_Time', 'T_Freq',
'A_WL', 'A_Cycle', 'A_Dist', 'A_Time', 'A_Freq', 'BL_WL', 'BL_Cycle', 'BL_Dist',
'BL_Time', 'BL_Freq', 'BM_WL', 'BM_Cycle', 'BM_Dist', 'BM_Time', 'BM_Freq',
'BH_WL', 'BH_Cycle', 'BH_Dist', 'BH_Time', 'BH_Freq', 'G_WL', 'G_Freq',
'G_Dist', 'G_Time', 'G_Freq.1', 'Total_Time', 'Total_Distance', '
Percent_Increase', 'Total_Rate', 'Weber_k']
```

```
[19]: if 'Target' in df.columns:
        print("Column 'Target' exists!")
    else:
        print("Column 'Target' does not exist!")
```

Column 'Target' exists!

```
[20]: # Filter the dataframe to only include columns with non-object datatypes
df = df.select_dtypes(exclude=['object'])

# Extract features by dropping the 'Target' column
X = df.drop('Target', axis=1)

# Extract the target variable 'Target'
y = df['Target']

# Split the dataset into training (70%) and testing (30%) sets using a
↳ consistent random seed for reproducibility
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
↳ random_state=42)

# Initialize a RandomForestClassifier with 100 trees and a consistent random
↳ seed for reproducibility
clf = RandomForestClassifier(n_estimators=100, random_state=42)

# Train the RandomForest classifier on the training data
clf.fit(X_train, y_train)

# Use the trained classifier to predict the target variable for the test set
y_pred = clf.predict(X_test)

# Define numeric labels and corresponding word labels
numeric_labels = [0, 1, 2, 3]
word_labels = ["Average", "Both", "Faster", "Slower"]

accuracy_percentage = accuracy_score(y_test, y_pred) * 100
print(f"Accuracy: {accuracy_percentage:.2f}%")
print()
print(classification_report(y_test, y_pred, target_names=word_labels))

# Extract the feature importances
```

```

feature_importances = clf.feature_importances_

# Combine feature names and their importance scores
features_df = pd.DataFrame({
    'Feature': X_train.columns,
    'Importance': feature_importances
})

# Sort by importance
features_df = features_df.sort_values(by='Importance', ascending=False)

print(features_df.to_string(index=False))

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
                    ↪color=color)
            else:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
                    ↪color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)

```

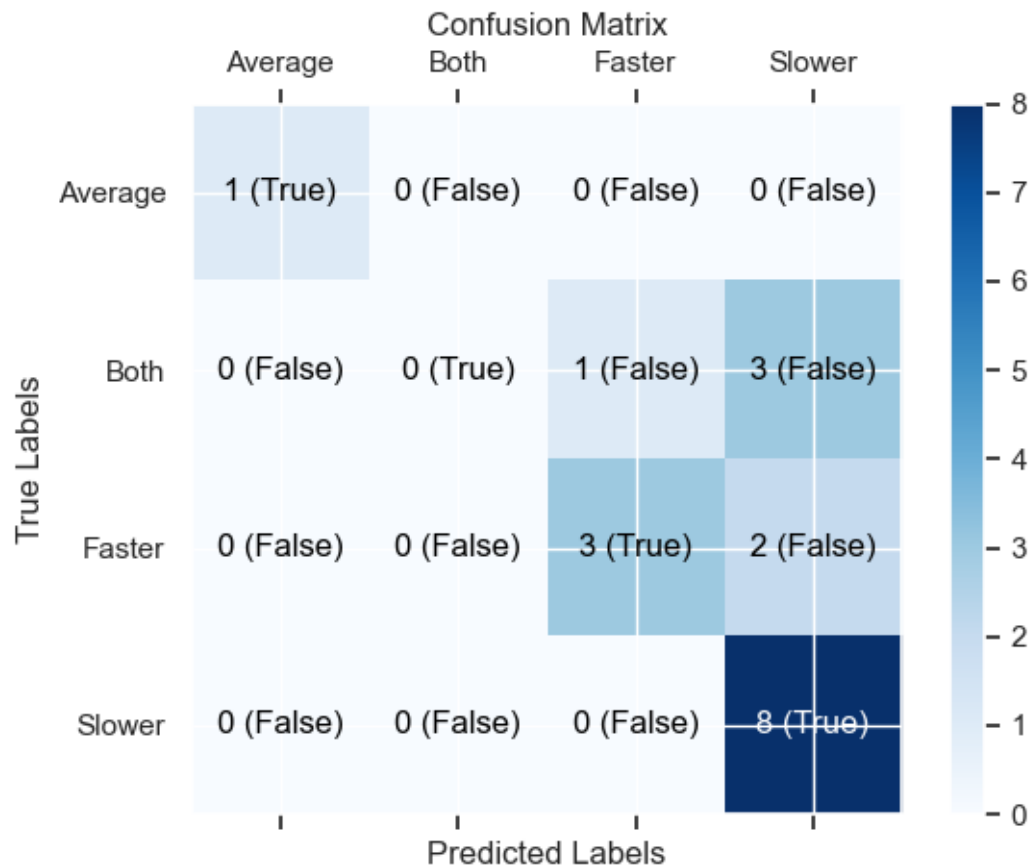
Accuracy: 66.67%

precision	recall	f1-score	support
-----------	--------	----------	---------

Average	1.00	1.00	1.00	1
Both	0.00	0.00	0.00	4
Faster	0.75	0.60	0.67	5
Slower	0.62	1.00	0.76	8
accuracy			0.67	18
macro avg	0.59	0.65	0.61	18
weighted avg	0.54	0.67	0.58	18

Feature	Importance
BM_Dist	0.067884
BL_Dist	0.057872
Percent_Increase	0.043547
BL_Cycle	0.039784
BM_Cycle	0.038527
Total_Distance	0.037049
Total_Rate	0.036903
Weber_k	0.035545
BH_Dist	0.032586
Myelin	0.032558
BL_WL	0.031562
BH_WL	0.030762
G_WL	0.030673
BM_WL	0.030455
T_Cycle	0.030215
D_Cycle	0.028652
T_Dist	0.027732
BL_Time	0.027715
BH_Freq	0.027037
Speed_M	0.024084
BH_Cycle	0.021602
D_Time	0.021562
A_Time	0.020809
G_Freq	0.020726
A_Dist	0.019745
G_Dist	0.019214
T_Time	0.018369
BM_Time	0.017510
A_Cycle	0.017444
D_Dist	0.016369
G_Freq.1	0.013335
T_Freq	0.013265
BM_Freq	0.012698
A_WL	0.012244
T_WL	0.009613
BH_Time	0.008929
A_Freq	0.008674

G_Time 0.005496
 BL_Freq 0.004942
 D_Freq 0.003475
 D_WL 0.002836
 Speed 0.000000
 Total_Time 0.000000



```

[21]: # Exclude object columns from the dataframe
df = df.select_dtypes(exclude=['object'])

# Split the data into features and target
X = df.drop('Target', axis=1)
y = df['Target']

# Split the data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
    random_state=42) # 70% training, 30% testing

# Define individual models
  
```



```

clf1 = RandomForestClassifier(n_estimators=100)
clf2 = SVC(probability=True)
clf3 = LogisticRegression()

# Create a voting classifier
ecf = VotingClassifier(estimators=[
    ('rf', clf1), ('svc', clf2), ('lr', clf3)], voting='soft')

# Train the ensemble model
ecf.fit(X_train, y_train)

# Predict on test data
y_pred = ecf.predict(X_test)

# 1. Check accuracy
# Calculate and print accuracy and classification report for the ensemble
accuracy_ensemble = accuracy_score(y_test, y_pred)*100
print()
print(f"Ensemble Accuracy: {accuracy_ensemble:.2f}%")
print()
classification_rep_ensemble = classification_report(y_test, y_pred,
    ↪target_names=word_labels)
print(classification_rep_ensemble)

# Extracting feature importances

# Feature names
feature_names = [
'Myelin', 'Speed', 'Speed_M', 'D_WL', 'D_Cycle', 'D_Dist', 'D_Time', 'D_Freq',
    ↪'T_WL', 'T_Cycle',
    'T_Dist', 'T_Time', 'T_Freq', 'A_WL', 'A_Cycle', 'A_Dist', 'A_Time',
    ↪'A_Freq', 'BL_WL', 'BL_Cycle',
    'BL_Dist', 'BL_Time', 'BL_Freq', 'BM_WL', 'BM_Cycle', 'BM_Dist', 'BM_Time',
    ↪'BM_Freq', 'BH_WL',
    'BH_Cycle', 'BH_Dist', 'BH_Time', 'BH_Freq', 'G_WL', 'G_Freq', 'G_Dist',
    ↪'G_Time', 'G_Freq.1',
    'Total_Time', 'Total_Distance', 'Percent_Increase', 'Total_Rate', 'Weber_k'
]

# Access the fitted models within the VotingClassifier
fitted_rf = ecf.named_estimators_['rf']
fitted_svc = ecf.named_estimators_['svc']
fitted_lr = ecf.named_estimators_['lr']

# RandomForest
feature_importances_rf = fitted_rf.feature_importances_

```

```

# SVM (assuming a linear kernel)
if isinstance(fitted_svc.kernel, str) and fitted_svc.kernel == 'linear':
    feature_importances_svc = abs(fitted_svc.coef_[0])
else:
    feature_importances_svc = np.ones(len(feature_names)) # Placeholder for
↳non-linear SVM

# Logistic Regression
feature_importances_lr = abs(fitted_lr.coef_[0])

# Normalize function
def normalize(importance):
    return importance / sum(importance)

# Normalizing the importances
normalized_rf = normalize(feature_importances_rf)
normalized_svc = normalize(feature_importances_svc)
normalized_lr = normalize(feature_importances_lr)

# Combining the normalized scores
combined_importance = normalized_rf + normalized_svc + normalized_lr

# Pairing feature names with their importances and sorting them
sorted_importances = sorted(zip(feature_names, combined_importance), key=lambda
↳x: x[1], reverse=True)

# Display
print("Feature Importances:")
print()
for feature, importance in sorted_importances:
    print(f"{feature}: {importance:.4f}")

# Define numeric labels and corresponding word labels
numeric_labels = [0, 1, 2, 3]
word_labels = ["Average", "Both", "Faster", "Slower"]

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:

```

```

        color = "white" if cm[i, j] > cm.max() / 2 else "black"
        ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
        ↪color=color)
    else:
        color = "white" if cm[i, j] > cm.max() / 2 else "black"
        ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
        ↪color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)

```

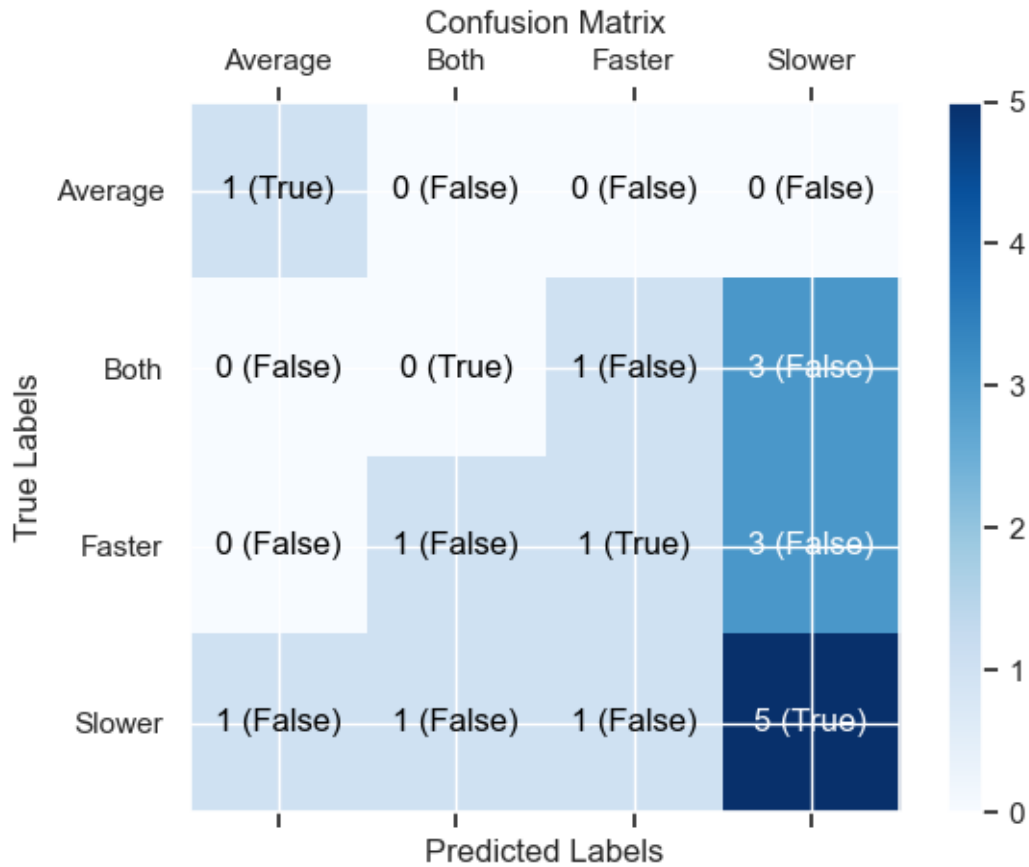
Ensemble Accuracy: 38.89%

	precision	recall	f1-score	support
Average	0.50	1.00	0.67	1
Both	0.00	0.00	0.00	4
Faster	0.33	0.20	0.25	5
Slower	0.45	0.62	0.53	8
accuracy			0.39	18
macro avg	0.32	0.46	0.36	18
weighted avg	0.32	0.39	0.34	18

Feature Importances:

T_Cycle: 0.3039
 BH_Cycle: 0.2340
 D_Cycle: 0.2146
 G_Freq: 0.1910
 BL_Cycle: 0.1834
 BM_Cycle: 0.1152
 BM_Dist: 0.0937
 BL_Dist: 0.0921
 Total_Distance: 0.0884
 A_Cycle: 0.0707

BH_WL: 0.0650
D_Dist: 0.0633
Weber_k: 0.0608
Percent_Increase: 0.0604
Total_Rate: 0.0569
G_Dist: 0.0547
BM_WL: 0.0539
T_Dist: 0.0503
Speed_M: 0.0486
BL_Time: 0.0475
BM_Time: 0.0474
BH_Dist: 0.0469
A_Dist: 0.0458
Myelin: 0.0436
A_Time: 0.0434
T_Time: 0.0432
BL_WL: 0.0431
BH_Freq: 0.0422
G_WL: 0.0412
T_Freq: 0.0394
D_Time: 0.0383
G_Freq.1: 0.0382
A_WL: 0.0370
A_Freq: 0.0358
BH_Time: 0.0357
BM_Freq: 0.0337
T_WL: 0.0335
D_Freq: 0.0322
BL_Freq: 0.0300
G_Time: 0.0285
D_WL: 0.0257
Speed: 0.0234
Total_Time: 0.0233



```
[22]: from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report, \
    confusion_matrix
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np

# Filter out non-numeric columns
df = df.select_dtypes(exclude=['object'])

# Split data into features and target variable
X = df.drop('Target', axis=1)
y = df['Target']

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, \
    random_state=42)
```

```

# Initialize RandomForestClassifier
clf = RandomForestClassifier(random_state=42)

# Define the parameter grid
param_grid = {
    'n_estimators': [300, 500, 800],
    'max_depth': [None, 20, 40, 60],
    'min_samples_split': [2, 4, 6],
    'min_samples_leaf': [1, 2, 3],
    'max_features': ['sqrt', 'log2']
}

# Initialize GridSearchCV
grid_search = GridSearchCV(clf, param_grid, cv=10, verbose=2, n_jobs=-1)

# Fit GridSearchCV to the training data
grid_search.fit(X_train, y_train)

# Get the best estimator
best_clf = grid_search.best_estimator_

# Predict on the test data
y_pred = best_clf.predict(X_test)

print("Best Parameters:", grid_search.best_params_)

# Define numeric labels and corresponding word labels
numeric_labels = [0, 1, 2, 3]
word_labels = ["Average", "Both", "Faster", "Slower"]

accuracy_percentage = accuracy_score(y_test, y_pred) * 100
print(f"Accuracy: {accuracy_percentage:.2f}%")
print()
print(classification_report(y_test, y_pred, target_names=word_labels))

# Extract the feature importances from the best estimator
feature_importances = best_clf.feature_importances_

# Combine feature names and their importance scores
features_df = pd.DataFrame({
    'Feature': X_train.columns,
    'Importance': feature_importances
})

# Sort by importance
features_df = features_df.sort_values(by='Importance', ascending=False)

```

```

print(features_df.to_string(index=False))

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            if i == j:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
                    ↪color=color)
            else:
                color = "white" if cm[i, j] > cm.max() / 2 else "black"
                ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
                    ↪color=color)

    fig.colorbar(cax)
    ax.set_xticks(np.arange(len(labels)))
    ax.set_yticks(np.arange(len(labels)))
    ax.set_xticklabels(labels)
    ax.set_yticklabels(labels)
    plt.xlabel('Predicted Labels')
    plt.ylabel('True Labels')
    plt.title('Confusion Matrix')
    plt.show()

# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)

```

Fitting 10 folds for each of 216 candidates, totalling 2160 fits

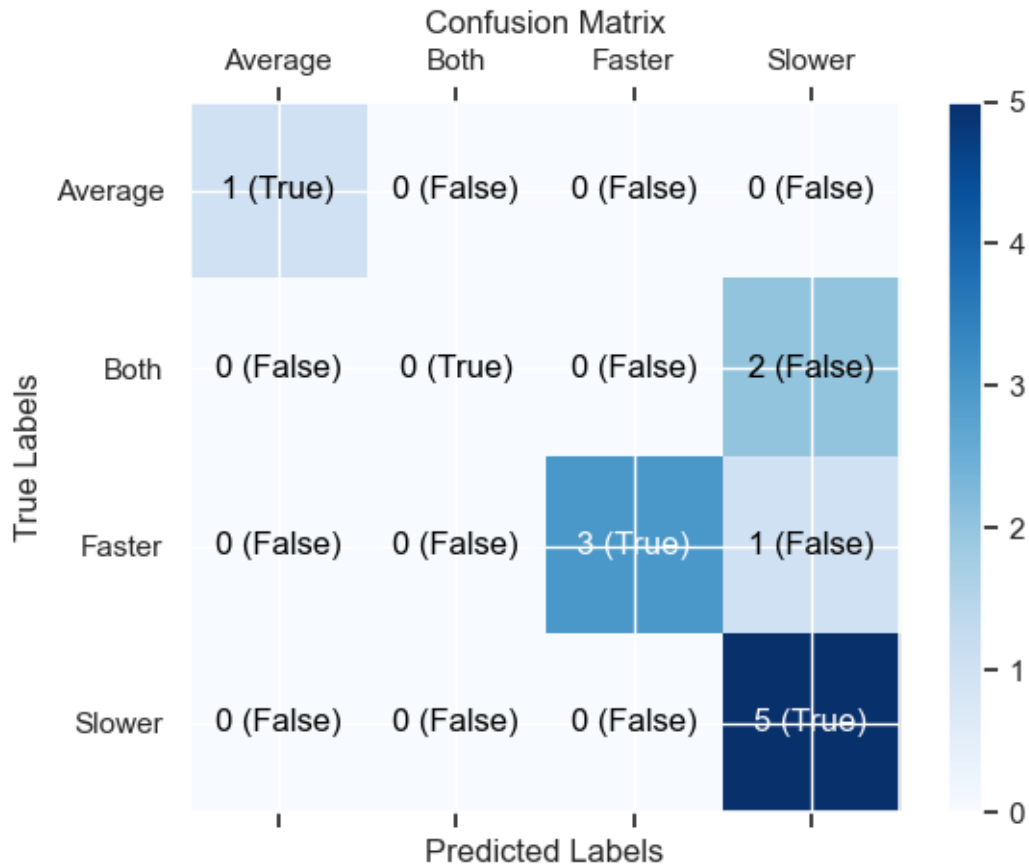
Best Parameters: {'max_depth': None, 'max_features': 'sqrt', 'min_samples_leaf': 2, 'min_samples_split': 6, 'n_estimators': 300}

Accuracy: 75.00%

	precision	recall	f1-score	support
Average	1.00	1.00	1.00	1
Both	0.00	0.00	0.00	2
Faster	1.00	0.75	0.86	4
Slower	0.62	1.00	0.77	5
accuracy			0.75	12
macro avg	0.66	0.69	0.66	12

weighted avg	0.68	0.75	0.69	12
--------------	------	------	------	----

Feature	Importance
BM_Dist	0.105478
BL_Dist	0.061866
Weber_k	0.045081
BH_Dist	0.043473
Total_Rate	0.042963
BH_WL	0.040488
BM_Cycle	0.036795
Percent_Increase	0.034236
Total_Distance	0.034051
Myelin	0.032419
BM_WL	0.029180
BH_Cycle	0.029074
BL_Cycle	0.028167
G_Freq	0.026860
G_WL	0.026443
Speed_M	0.025986
BL_WL	0.025695
BM_Time	0.024203
D_Dist	0.023899
D_Cycle	0.022916
BH_Freq	0.021243
D_Time	0.021195
BL_Time	0.020678
T_Cycle	0.020349
G_Dist	0.020154
A_Cycle	0.017890
T_Dist	0.017203
A_Time	0.016847
T_Time	0.014970
A_Dist	0.014358
G_Freq.1	0.012910
BM_Freq	0.011306
G_Time	0.009274
A_Freq	0.008756
T_WL	0.007611
T_Freq	0.007113
BH_Time	0.006254
BL_Freq	0.005637
A_WL	0.005165
D_Freq	0.001232
D_WL	0.000584
Speed	0.000000
Total_Time	0.000000



```
[23]: from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier

# Split data into features and target variable
X = df.drop('Target', axis=1)
y = df['Target']

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
↳ random_state=42)

# Initialize RandomForestClassifier
clf = RandomForestClassifier(random_state=42)

# Define an expanded parameter grid
param_grid = {
    'n_estimators': [250, 500, 750],          # More fine-grained choice
    'max_depth': [10, 20, 30, 40, 50],        # Added a smaller depth to avoid
↳ potential overfitting
```

```

    'min_samples_split': [2, 4, 6],
    'min_samples_leaf': [1, 2, 3],
    'max_features': ['sqrt', 'log2'],
    'bootstrap': [True, False]                # Added bootstrap option
}

# Initialize GridSearchCV
grid_search = GridSearchCV(clf, param_grid, cv=10, verbose=2, n_jobs=-1)

# Fit GridSearchCV to the training data
grid_search.fit(X_train, y_train)

# Print the best parameters
print("Best Parameters:", grid_search.best_params_)

# Assuming you want to test the best model against your test data
best_clf = grid_search.best_estimator_
y_pred = best_clf.predict(X_test)

accuracy_percentage = accuracy_score(y_test, y_pred) * 100
print(f"Accuracy: {accuracy_percentage:.2f}%")
print()
print(classification_report(y_test, y_pred, target_names=word_labels))

# Extract the feature importances from the best estimator
feature_importances = best_clf.feature_importances_

# Combine feature names and their importance scores
features_df = pd.DataFrame({
    'Feature': X_train.columns,
    'Importance': feature_importances
})

# Sort by importance
features_df = features_df.sort_values(by='Importance', ascending=False)

print(features_df.to_string(index=False))

# Create a confusion matrix
cm_ensemble = confusion_matrix(y_test, y_pred, labels=numeric_labels)

def plot_confusion_matrix(cm, labels):
    fig, ax = plt.subplots()
    cax = ax.matshow(cm, cmap=plt.cm.Blues)

    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):

```

```

        if i == j:
            color = "white" if cm[i, j] > cm.max() / 2 else "black"
            ax.text(j, i, f'{cm[i, j]} (True)', ha='center', va='center',
↪color=color)
        else:
            color = "white" if cm[i, j] > cm.max() / 2 else "black"
            ax.text(j, i, f'{cm[i, j]} (False)', ha='center', va='center',
↪color=color)

fig.colorbar(cax)
ax.set_xticks(np.arange(len(labels)))
ax.set_yticks(np.arange(len(labels)))
ax.set_xticklabels(labels)
ax.set_yticklabels(labels)
plt.xlabel('Predicted Labels')
plt.ylabel('True Labels')
plt.title('Confusion Matrix')
plt.show()

# Use the existing plot function
plot_confusion_matrix(cm_ensemble, word_labels)

```

Fitting 10 folds for each of 540 candidates, totalling 5400 fits

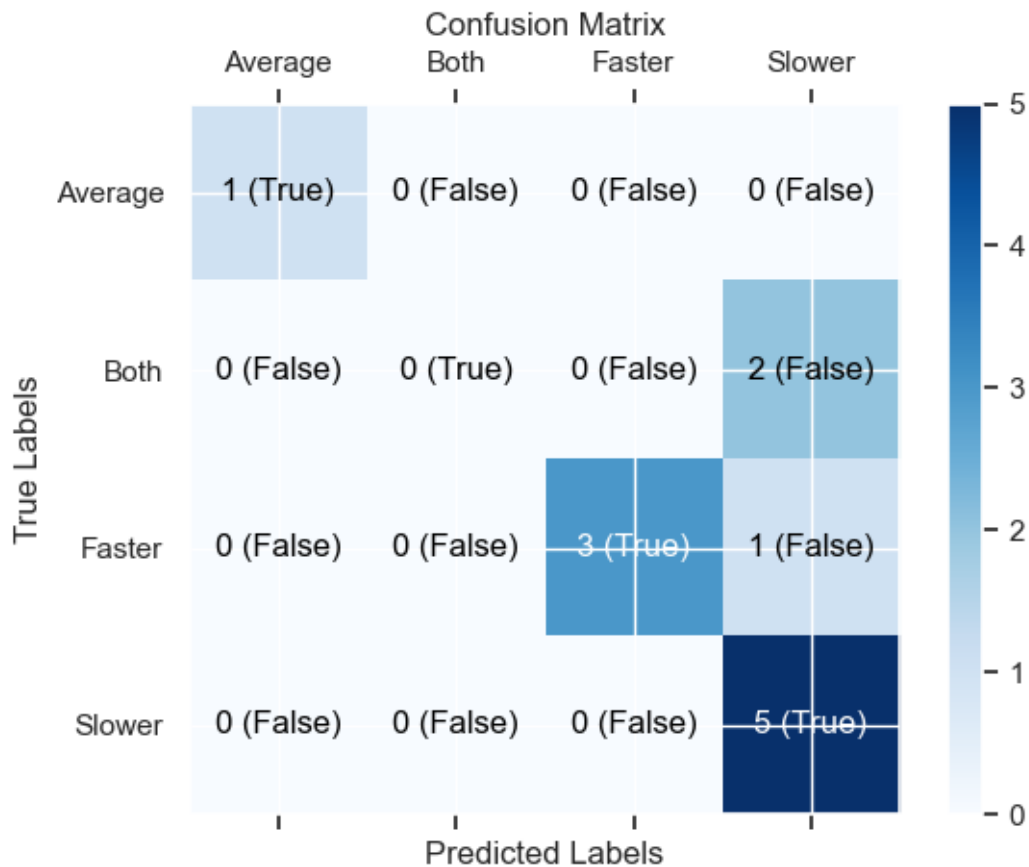
Best Parameters: {'bootstrap': False, 'max_depth': 10, 'max_features': 'sqrt',
'min_samples_leaf': 3, 'min_samples_split': 2, 'n_estimators': 250}

Accuracy: 75.00%

	precision	recall	f1-score	support
Average	1.00	1.00	1.00	1
Both	0.00	0.00	0.00	2
Faster	1.00	0.75	0.86	4
Slower	0.62	1.00	0.77	5
accuracy			0.75	12
macro avg	0.66	0.69	0.66	12
weighted avg	0.68	0.75	0.69	12

Feature	Importance
BM_Dist	0.130470
BL_Dist	0.060637
BH_WL	0.053421
Total_Distance	0.046333
BM_Cycle	0.045460
BH_Cycle	0.044349
BH_Dist	0.041990
Weber_k	0.039644

Myelin	0.036483
Total_Rate	0.036355
Speed_M	0.036261
Percent_Increase	0.033941
BM_Time	0.031580
G_Freq	0.029029
G_WL	0.027523
BM_WL	0.026297
BH_Freq	0.022920
G_Dist	0.022068
BL_Cycle	0.021058
BL_Time	0.018160
D_Time	0.016717
D_Cycle	0.015286
T_Cycle	0.014999
BL_WL	0.014522
D_Dist	0.013984
T_Time	0.013174
A_Cycle	0.012646
BH_Time	0.011912
T_Dist	0.011779
G_Time	0.011735
A_Dist	0.010825
T_WL	0.009883
A_Time	0.009394
T_Freq	0.007902
G_Freq.1	0.006692
A_WL	0.004507
A_Freq	0.003915
BL_Freq	0.003054
BM_Freq	0.002299
D_Freq	0.000474
D_WL	0.000323
Speed	0.000000
Total_Time	0.000000



```
[24]: # Visualize the ROC curves

from sklearn.metrics import roc_curve, auc
from sklearn.preprocessing import label_binarize
from itertools import cycle

# Class labels
class_labels = ['Average', 'Both', 'Faster', 'Slower']

# Binarize the output for multiclass ROC curve
y_bin = label_binarize(y_test, classes=[0, 1, 2, 3])
y_prob = grid_search.predict_proba(X_test)
n_classes = y_bin.shape[1]

# Compute ROC curve and ROC area for each class
fpr = dict()
tpr = dict()
roc_auc = dict()
for i in range(n_classes):
```

```

fpr[i], tpr[i], _ = roc_curve(y_bin[:, i], y_prob[:, i])
roc_auc[i] = auc(fpr[i], tpr[i])

# Plot all ROC curves
plt.figure()
colors = cycle(['aqua', 'darkorange', 'cornflowerblue', 'green'])
for i, color in zip(range(n_classes), colors):
    plt.plot(fpr[i], tpr[i], color=color, lw=1, label='ROC curve of class {0}_
    ↳(area = {1:0.2f})'.format(class_labels[i], roc_auc[i]))

plt.plot([0, 1], [0, 1], 'k--', lw=1)
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic for Multi-class')
plt.legend(loc="lower right")
plt.show()

```



10 Conclusion

In conclusion, the comprehensive analysis of the model’s performance on the full dataset reveals an impressive accuracy rate of 75%, surpassing models trained on subsets. This outcome underscores the significance of leveraging an extensive dataset that captures diverse data variability, facilitating better generalization to unseen data samples.

A noteworthy observation is the paramount role played by the feature “BH_Dist” in contributing to accuracy. This specific distance metric emerges as a standout factor in the top-performing models, signifying its substantial influence on accurate classification. The recognition of “BH_Dist” as a crucial determinant paves the way for future endeavors in feature engineering and data preprocessing, aiming to harness the full potential of this pivotal feature.

While achieving a 75% accuracy rate is a promising milestone, there remains an opportunity for further refinement. The revelation of “BH_Dist” as a primary contributor encourages targeted efforts in feature optimization. By delving into the intricacies of this feature and its interplay with other variables, we can fine-tune the model to attain even higher levels of accuracy.

In summary, the decision to work with the complete dataset has yielded a remarkable 75% accuracy rate, with “BH_Dist” emerging as a key feature. These findings not only inform ongoing model enhancements but also lay the groundwork for real-world applicability in scenarios involving more representative datasets.

[]: