## Main.py

```python
1   import pandas as pd
2   import numpy as np
3
4   # Define paths for each dataset in your Google Drive
5   df1_path = '/Users/ankushbhatt/Downloads/VFS-CS/DataSet/groundtruth.csv'  # Update
    with your actual path
6   df2_path = '/Users/ankushbhatt/Downloads/VFS-CS/DataSet/weather_features.csv'  #
    Update with your actual path
7   df3_path = '/Users/ankushbhatt/Downloads/VFS-CS/DataSet/road_features.csv'  # Update
    with your actual path
8
9   # define the datatypes for each csv
10  dtypes_df1 = {'road_segment_id': 'int32', 'timestamp': 'str', 'max_capacity':
    'int32', 'occupied': 'int32', 'available': 'int32'}
11  dtypes_df2 = {'road_segment_id': 'int32', 'timestamp': 'str', 'tempC': 'float32',
    'windspeedKmph': 'float32', 'precipMM': 'float32'}
12  dtypes_df3 = {'road_segment_id': 'int32', 'commercial': 'Int64', 'residential':
    'Int64', 'transportation': 'Int64', 'schools': 'Int64', 'eventsites': 'Int64',
13              'restaurant': 'Int64', 'shopping': 'Int64', 'office': 'Int64',
    'supermarket': 'Int64', 'num_off_street_parking': 'Int64', 'off_street_capa':
    'Int64'}
14
15  # load the datasets using the defined paths
16  df1 = pd.read_csv(df1_path, dtype=dtypes_df1)
17  df2 = pd.read_csv(df2_path, dtype=dtypes_df2)
18  df3 = pd.read_csv(df3_path, dtype=dtypes_df3)
19
20
21  # merge the dataset
22  # first merge df1 and df2 on road_segment_id and timestamp
23  merged_df = pd.merge(df1, df2, on=['road_segment_id', 'timestamp'], how='inner')
24  # second merge df3 with merged_df on road_segment_id
25  final_df = pd.merge(merged_df, df3, on='road_segment_id', how='inner')
26
27  # save the merged dataset
28  final_df.to_csv('New_joint_dataset.csv', index=False)
29
30  # Print the final merged dataset info
31  print("\nMerged Dataset:")
32  print(final_df.info()) # showing information about final merged dataset
33  print(final_df.head()) # showing first few roows of finam merged dataset
34
35  # Data Pre processing
36
37  # Load New Joint dataset
38
39  import pandas as pd
40  import numpy as np
41  df = pd.read_csv('New_joint_dataset.csv')
42
43  # get a summary of the data
44  print("First few rows of dataframe")
```

```python
45 │ print(df.head())
46 │ print("\ndataframe info:")
47 │ print(df.info())
48 │ print("\nsummary of dataframe:")
49 │ print(df.describe(include='all'))
50 │
51 │ # Check for missing values
52 │ missing_values = df.isnull().sum()
53 │ print("\nMissing Values:")
54 │ print(missing_values)
55 │
56 │ # Fill the missing values
57 │ # Now fill numeric values with median and categorical with mode
58 │ for column in final_df.columns:
59 │   if df[column].dtype in ['float64', 'int64']:
60 │     df[column].fillna(df[column].median(), inplace=True)
61 │   elif df[column].dtype == 'object':
62 │     df[column].fillna(df[column].mode()[0], inplace=True)
63 │
64 │ # Identify and handle the outliers
65 │ # Now using IQR methods for detecting outliers in the data
66 │ def remove_outliers_iqr(data):
67 │   for column in data.select_dtypes(include=['float64', 'int64']).columns:
68 │     Q1 = data[column].quantile(0.25)
69 │     Q3 = data[column].quantile(0.75)
70 │     IQR = Q3 - Q1
71 │     lower_bound = Q1 - 1.5 * IQR
72 │     upper_bound = Q3 + 1.5 * IQR
73 │     data = data[(data[column] >= lower_bound) & (data[column] <= upper_bound)]
74 │   return data
75 │ df = remove_outliers_iqr(df)
76 │
77 │ # Standardize the format of data
78 │ df['timestamp'] = pd.to_datetime(df['timestamp'], errors='coerce')  # 'coerce' will
   │ turn invalid parsing into NaT (Not a Time)
79 │
80 │ # Extract useful date-time features
81 │ df['year'] = df['timestamp'].dt.year
82 │ df['month'] = df['timestamp'].dt.month
83 │ df['day'] = df['timestamp'].dt.day
84 │ df['hour'] = df['timestamp'].dt.hour
85 │ df['minute'] = df['timestamp'].dt.minute
86 │ df['dayofweek'] = df['timestamp'].dt.dayofweek  # Monday=0, Sunday=6
87 │
88 │ # Drop the original 'timestamp' column as it is no longer needed for prediction
89 │ df.drop(columns=['timestamp'], inplace=True)
90 │
91 │ # Now you can proceed with other steps (e.g., missing values handling, outliers
   │ removal, etc.)
92 │
93 │ # Now removing duplicates
94 │ df.drop_duplicates(inplace=True)
95 │
96 │ # dropping the unnecessary columns
```

```python
97   df.drop(columns=['Unnamed: 0_x', 'Unnamed: 0_y', 'Unnamed: 0', 'max_capacity',
     'occupied'], inplace=True, errors='ignore')
98   print("Columns after dropping:", df.columns)
99
100  # convert data types
101  # convert numeric columns to appropriate types
102  from sklearn.preprocessing import StandardScaler, PowerTransformer
103  numeric_columns = ['max_capacity', 'occupied', 'available', 'tempC',
     'windspeedKmph', 'precipMM', 'commercial', 'residential', 'transportation',
     'schools', 'eventsites', 'restaurant', 'shopping',
104                      'office', 'supermarket', 'num_off_street_parking',
     'off_street_capa']
105  for column in numeric_columns:
106    if column in df.columns:
107      df[column] = pd.to_numeric(df[column], errors='coerce')
108
109  # convert 'road_segment_id' to int32 if necesaary
110  if 'road_segment_id' in df.columns:
111    df['road_segment_id'] = df['road_segment_id'].astype('int32')
112
113  #**Skewness Correction with Power Transformation**
114  # Apply PowerTransformer to correct skewness in numeric columns
115  numeric_features = df.select_dtypes(include=['float64', 'int64']).columns
116  power_transformer = PowerTransformer(method='yeo-johnson', standardize=True)
117  df[numeric_features] = power_transformer.fit_transform(df[numeric_features])
118
119  # **Scaling Features**
120  # Use StandardScaler to standardize numeric features
121  scaler = StandardScaler()
122  df[numeric_features] = scaler.fit_transform(df[numeric_features])
123
124  # check datatypes after conversion
125  print("\nData Types after Conversion:")
126  print(df.dtypes)
127
128  # check the info of the final dataframe
129  print("\nFinal Dataframe Info:")
130  print(df.info())
131
132  # save the clean data
133  clean_data_path = 'clean_data.csv'
134  df.to_csv(clean_data_path, index=False)
135  print("\nData preprocessing completed. clean data saved to", clean_data_path)
136
137  # Trained the Model:
138  import pandas as pd
139  import numpy as np
140  from sklearn.model_selection import train_test_split, RandomizedSearchCV,
     cross_val_score
141  from sklearn.ensemble import RandomForestRegressor
142  from sklearn.metrics import mean_absolute_error, mean_squared_error
143  import matplotlib.pyplot as plt
144
145  # Load the cleaned dataset
```

```python
146  cd = pd.read_csv('clean_data.csv')
147
148  # Features and Target
149  x = cd.drop(columns=['available'])  # Features
150  y = cd['available']  # Target
151
152  # Split the data into train and test sets (80% train, 20% test)
153  x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2,
     random_state=42)
154
155  # Define the model with slightly fewer trees and shallower depth
156  model = RandomForestRegressor(
157      n_estimators=30,            # Reduced n_estimators for faster learning
158      max_depth=5,                # Limiting tree depth to avoid overfitting
159      min_samples_split=5,        # Minimum samples required to split a node
160      min_samples_leaf=4,         # Minimum samples required to be in a leaf node
161      random_state=42,
162      n_jobs=-1                    # Use all available CPU cores for faster training
163  )
164
165  # Fit the model on training data
166  model.fit(x_train, y_train)
167
168  # Predict on test data
169  y_pred = model.predict(x_test)
170
171  # Evaluate the model
172  mae = mean_absolute_error(y_test, y_pred)
173  rmse = np.sqrt(mean_squared_error(y_test, y_pred))
174  r2_score = model.score(x_test, y_test)
175
176  # Print evaluation metrics
177  print("Model Evaluation Metrics:")
178  print("Mean Absolute Error (MAE):", mae)
179  print("Root Mean Squared Error (RMSE):", rmse)
180  print("R² (Coefficient of Determination) accuracy:", r2_score)
181
182  # Cross-validation with fewer folds (3-fold to save time)
183  cv_scores = cross_val_score(model, x, y, cv=5, scoring='neg_mean_squared_error',
     n_jobs=-1)
184  print("\nCross-validation scores:", cv_scores)
185  print("Average Cross-validation score:", cv_scores.mean())
186
187  # Hyperparameter tuning with a focused parameter grid and reduced combinations
188  param_distributions = {
189      'n_estimators': [10, 30, 50],  # Reduced range for faster tuning
190      'max_depth': [5, 7],           # Limiting to fewer options
191      'min_samples_split': [2, 5],   # Narrowing range
192      'min_samples_leaf': [1, 2, 4], # Fewer leaf node options
193  }
194
195  # RandomizedSearchCV with fewer iterations
196  randomized_search = RandomizedSearchCV(
197      RandomForestRegressor(random_state=42, n_jobs=-1),
```

```python
198        param_distributions,
199        n_iter=5,                          # Fewer combinations to try
200        cv=5,                              # Reduced folds
201        scoring='neg_mean_squared_error',
202        n_jobs=-1
203    )
204
205    # Fit RandomizedSearchCV to training data
206    randomized_search.fit(x_train, y_train)
207
208    # Print best hyperparameters from RandomizedSearchCV
209    print("Best Hyperparameters from RandomizedSearchCV:")
210    print(randomized_search.best_params_)
211
212    # Retrain the model with the best parameters
213    best_model = randomized_search.best_estimator_
214    best_model.fit(x_train, y_train)
215
216    # Predict and evaluate the retrained model
217    y_pred_best = best_model.predict(x_test)
218
219    # Evaluate the model
220    mae_best = mean_absolute_error(y_test, y_pred_best)
221    rmse_best = np.sqrt(mean_squared_error(y_test, y_pred_best))
222    r2_score_best = best_model.score(x_test, y_test)
223
224    # Print evaluation metrics of the best model
225    print("\nBest Model Evaluation Metrics:")
226    print("Mean Absolute Error (MAE):", mae_best)
227    print("Root Mean Squared Error (RMSE):", rmse_best)
228    print("R² (Coefficient of Determination) accuracy:", r2_score_best)
229
230    # Final model accuracy
231    accuracy_percentage_best = r2_score_best * 100
232    print(f"Final Accuracy: {accuracy_percentage_best:.2f}%")
233
234    # Cross-validation to evaluate model performance on multiple folds (using 3-fold CV)
235    cv_scores_best = cross_val_score(best_model, x, y, cv=5,
       scoring='neg_mean_squared_error', n_jobs=-1)
236    print("\nCross-validation scores of best model:", cv_scores_best)
237    print("Average Cross-validation score of best model:", cv_scores_best.mean())
238
```