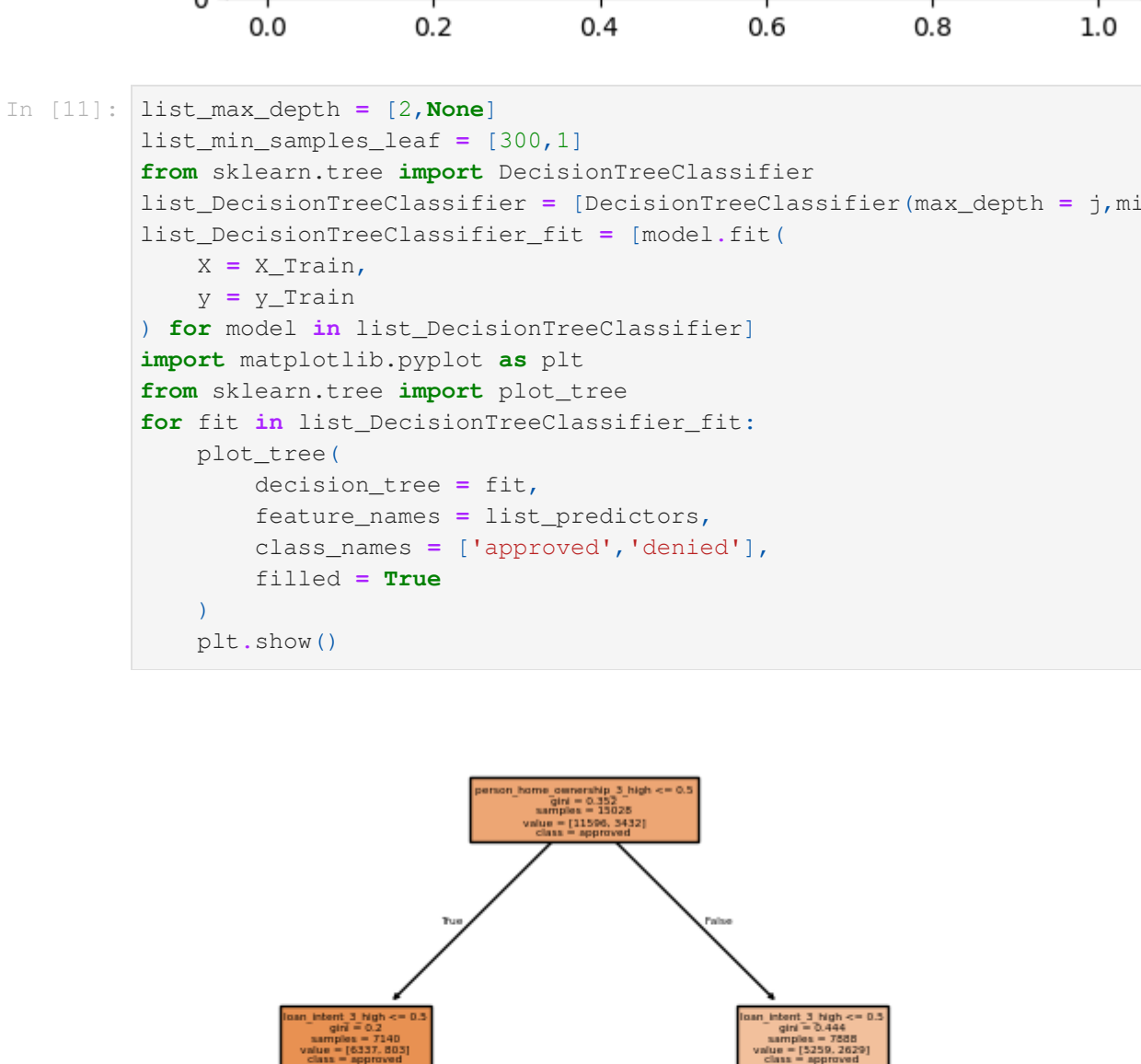


```

1 engine = 'pytorch'
2 df_clean.head()
3
4
5 In [10]:
6 loan_status, person_education_2_moderate, person_education_3_high, person_home_ownership_3_high, loan_intent_2_moderate, loan_intent_3_high
7
8 0 1 test 0 0 1 1 0 0 0
9
10 1 0 test 1 0 0 0 0 0 0
11
12 2 1 train 1 0 0 0 0 0 1
13
14 3 1 train 0 1 1 1 0 1 1
15
16 4 1 test 0 0 1 0 1 0 1
17
18
19 In [10]:
20 list_predictors = [
21     'person_education_2_moderate',
22     'person_education_3_high',
23     'person_home_ownership_3_high',
24     'loan_intent_2_moderate',
25     'loan_intent_3_high',
26 ]
27
28 X = df_clean[list_predictors]
29 y = df_clean['loan_status']
30 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
31 y_train = df_clean['loan_status'].iloc[0:df_clean['loan_status'].count()*0.8]
32 y_test = df_clean['loan_status'].iloc[df_clean['loan_status'].count()*0.8:]
33

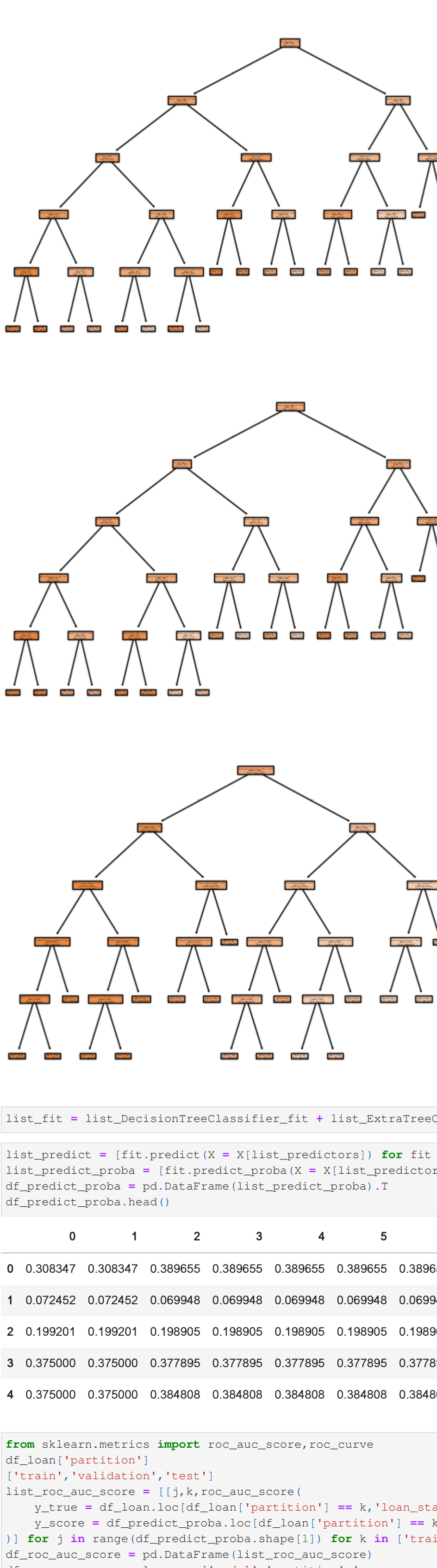
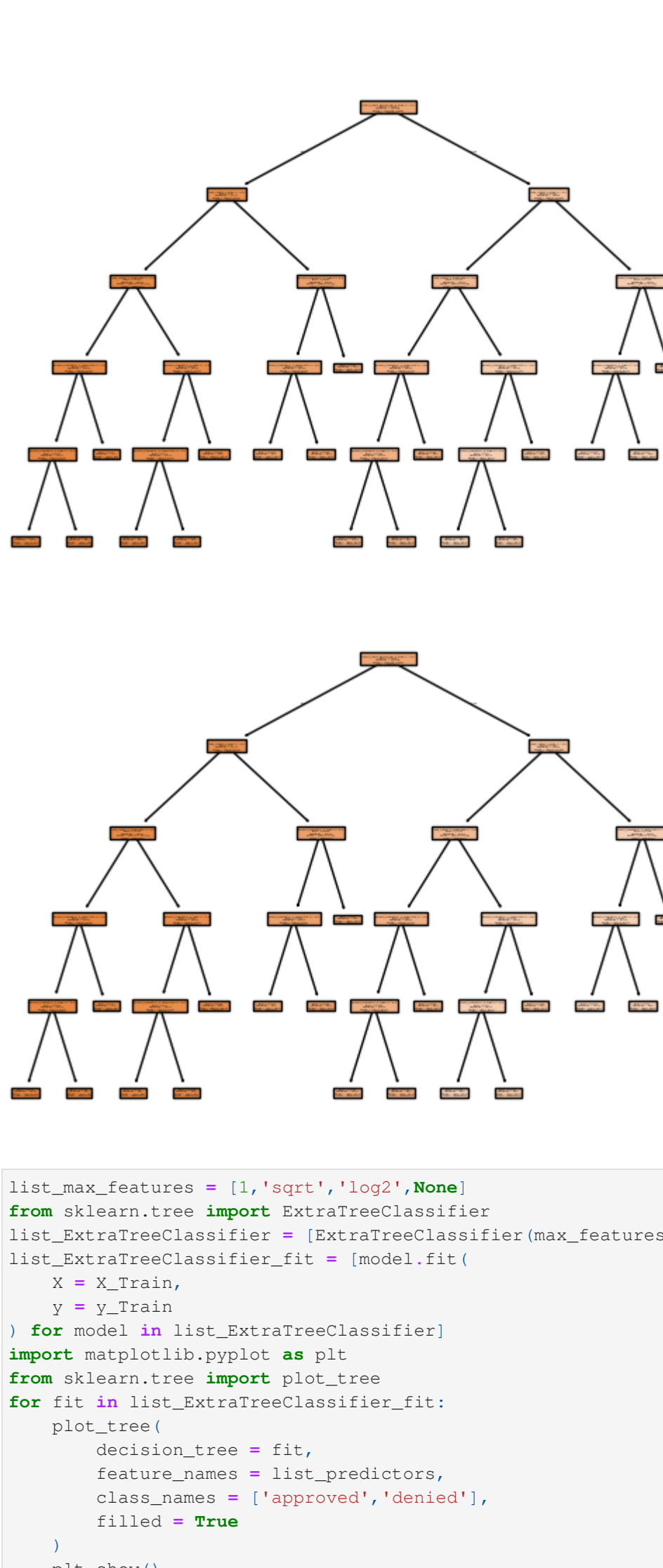
```

Category	Number of Publications
Other	~34,000
Non-peer-reviewed	~10,000



```

graph TD
    Root["number < 1000000"]
    Root -- Yes --> Leaf1["number < 1000000"]
    Root -- No --> Node1["number < 1000000"]
    Node1 -- Yes --> Leaf2["number < 1000000"]
    Node1 -- No --> Leaf3["number < 1000000"]
  
```



```

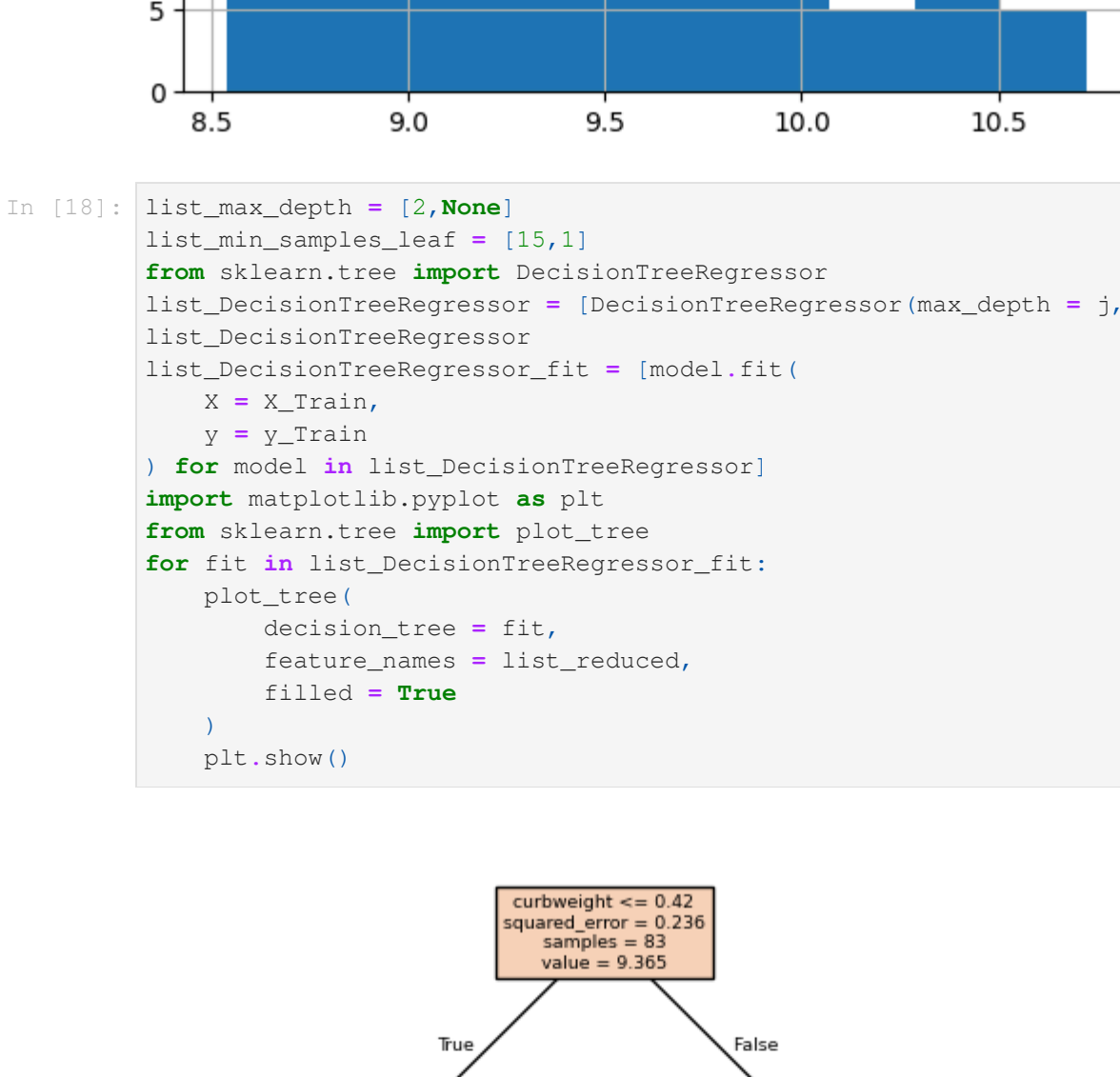
        column='partition',
        values='roc_auc_score'
    )
]

```

```
df_auc_auc_scoreort_values(['test','validation','train'])
```

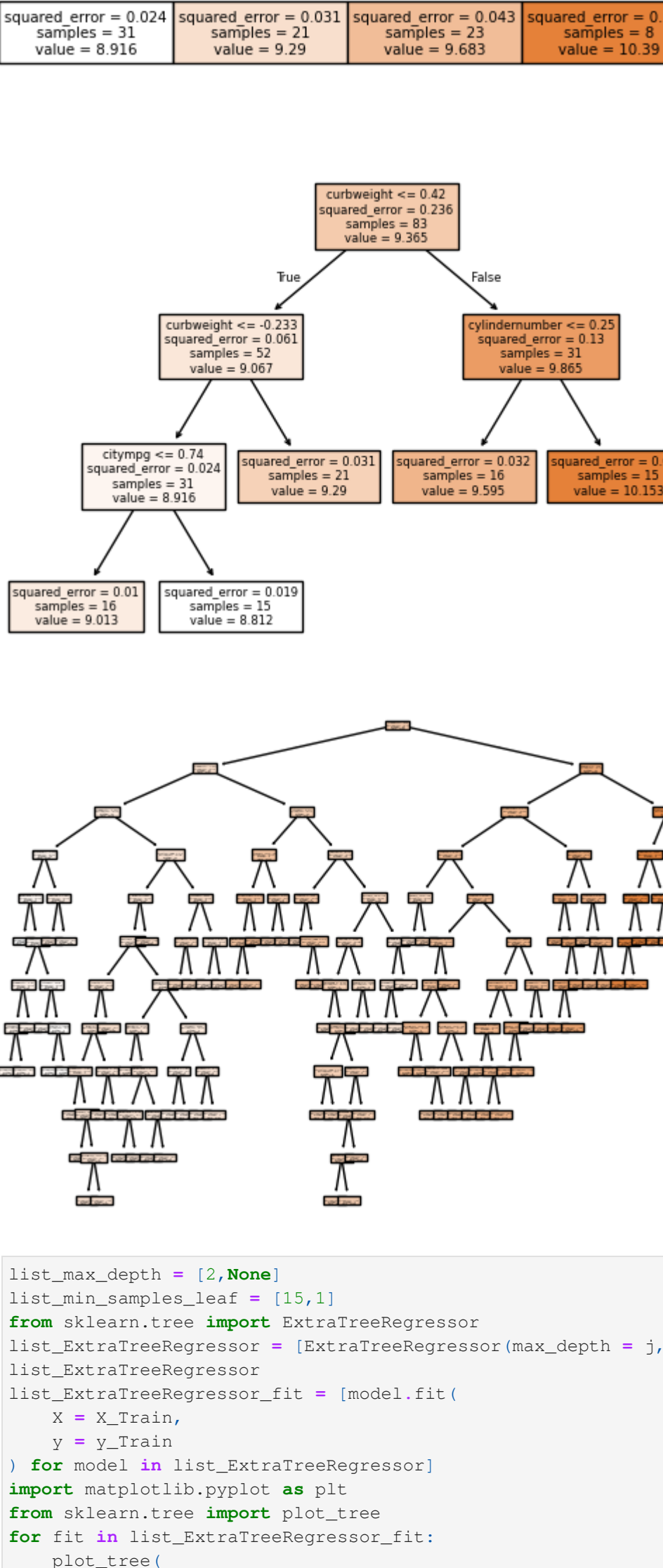
partition	test	train	validation	TreeClassifier
model				
0	0.681409	0.655993	0.679326	DecisionTreeClassifier(max_depth=2, min_sample...
1	0.681409	0.655993	0.679326	DecisionTreeClassifier(max_depth=2)

2	0.68807	0.699059	0.692372	DecisionTreeClassifier(min_samples_leaf=300)
3	0.68807	0.699059	0.692372	DecisionTreeClassifier()
4	0.68807	0.699059	0.692372	ExtraTreeClassifier(max_features=1)
5	0.68807	0.699059	0.692372	ExtraTreeClassifier()
6	0.68807	0.699059	0.692372	ExtraTreeClassifier(max_features=log2)
7	0.68807	0.699059	0.692372	ExtraTreeClassifier(max_features=None)

[illegible]

```

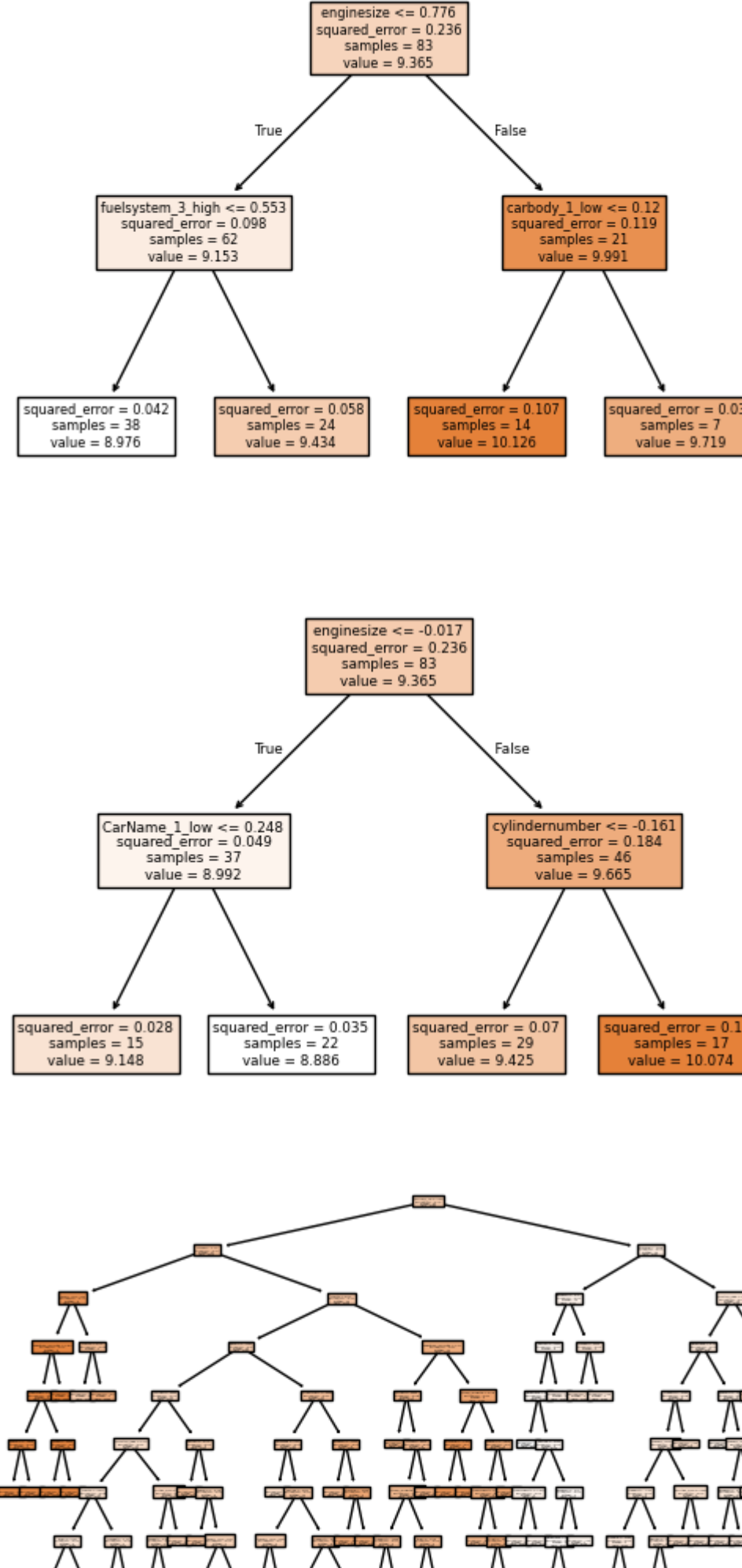
graph TD
    Root["speed = 3.067"]
    Root --> L1["speed_err = 0.024  
accident_prob = 0.925"]
    Root --> R1["speed_err = 0.031  
speed = 0.29"]
    R1 --> L2["speed_err = 0.035  
accident_prob = 0.935"]
    R1 --> R2["speed_err = 0.033  
accident_prob = 0.920"]
    Root --> R3["speed = 3.285"]
    R3 --> L3["speed_err = 0.1  
accident_prob = 0.935"]
    R3 --> R4["speed_err = 0.033  
accident_prob = 0.920"]
  
```



```

    }
    plt.show()

```



```
In [21]: list_predict = [fit.predict(X=X[list_reduced]) for fit in list_fit]
df_predict = pd.DataFrame(list_predict).T
df_predict.head()
```

	0	1	2	3	4	5	6	7
0	9.290163	9.290163	9.290163	9.510075	9.244562	9.433748	9.425308	9.510075
1	9.290163	9.290163	9.290163	9.510075	9.244562	9.433748	9.425308	9.510075
2	10.153301	9.682777	10.153301	9.999616	10.120741	9.719882	10.074106	9.632535

```
3  8.915628  8.915628  9.012965  9.206332  9.244562  9.433748  9.147842  9.296335
4  10.153301  9.682777  10.153301  8.999616  10.120741  9.433748  10.074106  9.632335
```

```
In [22]: from sklearn.metrics import mean_squared_error
```



