

```
In [7]: import numpy as np
import sklearn.datasets
import sklearn.linear_model
import sklearn.metrics

In [2]: np.random.seed(42) # do not change for reproducibility
```

## Task and problem statement

In this test we use the [Covertype Data Set](#), a dataset describing cartographic features of areas of land in the USA and also its forest type according to the US Forest Service. There are seven classes (1-7), 581012 samples and 54 features. For this test, we're only interested in cover type 3.

```
In [4]: dataset = sklearn.datasets.fetch_covtype()

In [6]: features = dataset.data
target = dataset.target
#feature_names = dataset.feature_names

In [7]: # only use a random subset for speed - pretend the rest of the data doesn't exist
random_sample = np.random.choice(len(dataset.data), len(dataset.data) // 10)

In [8]: len(random_sample)

Out[8]: 58101

In [48]: COVER_TYPE = 3
features = dataset.data[random_sample, :]
target = dataset.target[random_sample] == COVER_TYPE

In [49]: features.shape

Out[49]: (58101, 54)
```

A junior colleague tells you that they're getting 96% accuracy using logistic regression. You review their work and see the following:

```
In [8]: classifier = sklearn.linear_model.LogisticRegression(solver='liblinear') # do not change this

In [11]: classifier.fit(features, target)
preds = classifier.predict(features)
acc = sklearn.metrics.accuracy_score(preds, target)

In [12]: print(f'Accuracy: {acc:.3f}')
```

## Question 1

Comment on what was done above. Evaluate the accuracy more thoroughly. Do not modify the parameters of the model (the cell marked with the comment). Use the classifier object.

## Answer

In the above, the problem is transformed into a binary classification problem, thereby testing if an observation is either of cover type 3 or not. Then a logistic regression model is trained and evaluated, which achieved an accuracy of approximately 96%. However, the trained model is evaluated on the same data it was trained on, which reminds me of a funny ML meme !)

Further, the model was evaluated w.r.t accuracy metric. However, since the classes are imbalanced, using accuracy will give distorted performance view when we'll take a closer look for individual classes. Therefore, we need to interpret using balanced performance metrics such as f1, MCC, or AUC scores.

So, let's create separate training and test sets: 80:20 ratio. We train the model on training set and evaluate the model on held-out test set to see how accurately we can classify the test samples.

```
In [53]: from sklearn.model_selection import train_test_split

def getdataForCoverType3(dataset, random_sample, to_sampled=False):
    random_sample = np.random.choice(len(dataset.data), len(dataset.data) // 10)

    if to_sampled:
        features = dataset.data
        target = dataset.target == COVER_TYPE
    else:
        features = dataset.data[random_sample, :]
        target = dataset.target[random_sample] == COVER_TYPE

    return features, target

In [58]: COVER_TYPE = 3
features, target = getdataForCoverType3(dataset, COVER_TYPE, False)
X_train, X_test, y_train, y_test = train_test_split(features, target, test_size=0.2, random_state=42)
```

We train the model and generate predictions on unseen test set:

```
In [14]: classifier.fit(X_train, y_train)

I wrote a method, which generates different metrics to evaluate the predictive power of the model on unseen test set. I could include it as a util, but doing here for the time being.

In [16]: from sklearn.metrics import precision_score, recall_score, f1_score, roc_auc_score, accuracy_score, classification_report
from sklearn.metrics import plot_confusion_matrix, plot_roc_curve, plot_precision_recall_curve

def evaluateModel(classifier, X_test, y_test, weighted=False):
    """This method generates different metrics to evaluate the predictive power of the model on unseen test set"""
    y_pred = model.predict(X_test)

    plot_confusion_matrix(model, X_test, y_test)
    plt.show()

    plot_roc_curve(model, X_test, y_test)
    plt.show()

    plot_precision_recall_curve(model, X_test, y_test)
    plt.show()

    print(classification_report(y_test, y_pred))

    if weighted:
        acc = accuracy_score(y_test, y_pred)
        print(f'Accuracy: {acc:.3f}')

        precision = precision_score(y_test, y_pred, average='weighted')
        print(f'Precision: {precision:.3f}')

        recall = recall_score(y_test, y_pred, average='weighted')
        print(f'Recall: {recall:.3f}')

        f1 = f1_score(y_test, y_pred, average='weighted')
        print(f'F1 score: {f1:.3f}')

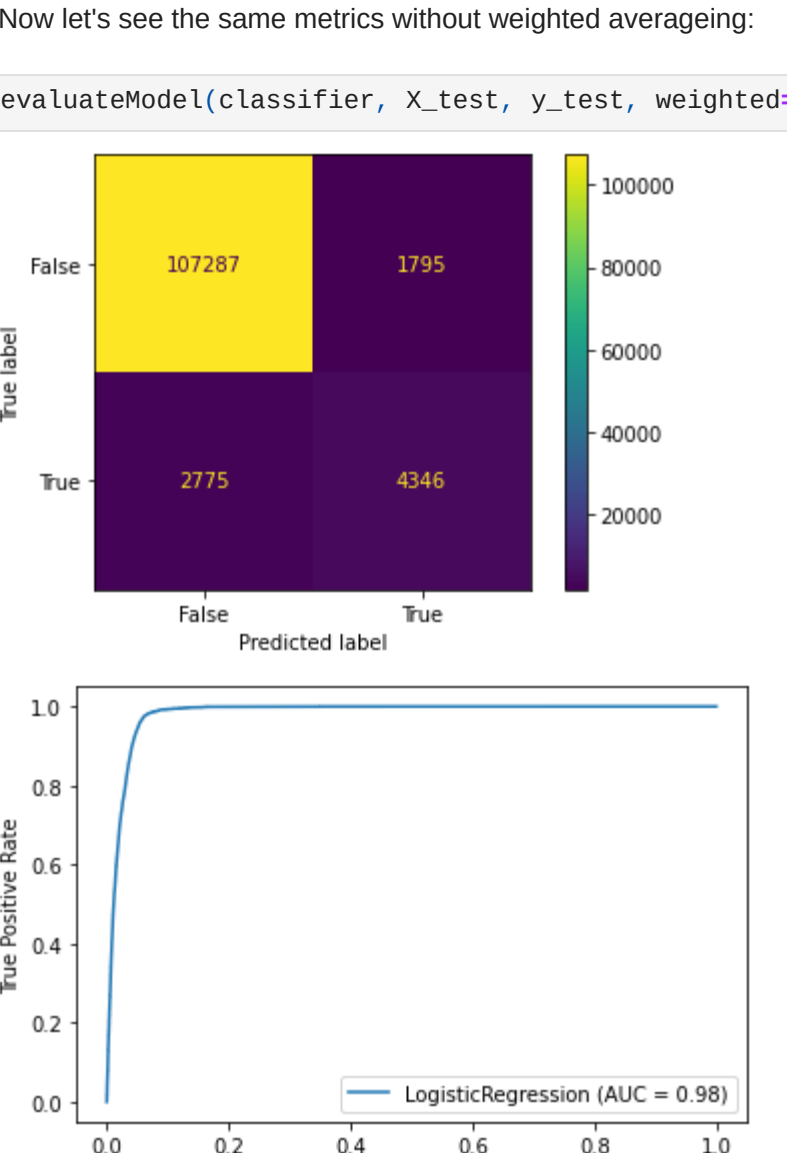
        auc = roc_auc_score(y_test, y_pred, average='weighted')
        print(f'AUC score: {auc:.3f}')

        mcc = matthews_corcoef(y_test, y_pred)
        print(f'MCC score: {mcc:.3f}')

    else:
        print(f'Accuracy: {accuracy_score(y_test, y_pred):.3f}')
        print(f'Precision: {precision_score(y_test, y_pred):.3f}')
        print(f'Recall: {recall_score(y_test, y_pred):.3f}')
        print(f'F1: {f1_score(y_test, y_pred):.3f}')
        print(f'AUC score: {roc_auc_score(y_test, y_pred):.3f}')
        print(f'MCC score: {matthews_corcoef(y_test, y_pred):.3f}')
```

Now if we take the weighted average scoring, we will see very different scenario:

```
In [39]: evaluateModel(classifier, X_test, y_test, weighted=True)
```



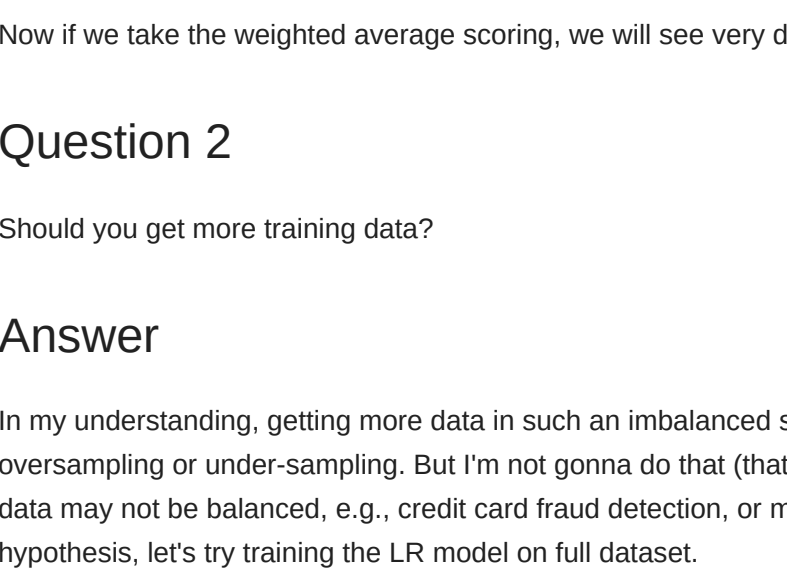
The confusion matrix shows 107287 True Positives, 1795 False Positives, 2775 False Negatives, and 4346 True Negatives. The ROC curve shows an AUC of 0.98, and the Precision-Recall curve shows an AP of 0.70. The classification report shows a weighted accuracy of 0.96, precision of 0.96, recall of 0.96, and F1 score of 0.96.

Thus having an weighted f1-score of 95.9% on test set is great. However, if we closely observe at the class-wise classification report, we can see that the classifier has made significant misclassification for the true class (cover type 3), making the scores much lower compared than that of false class. This is obvious because classes are severely imbalanced.

```
In [40]: import numpy as np

unique, counts = np.unique(target, return_counts=True)
dict(zip(unique, counts))

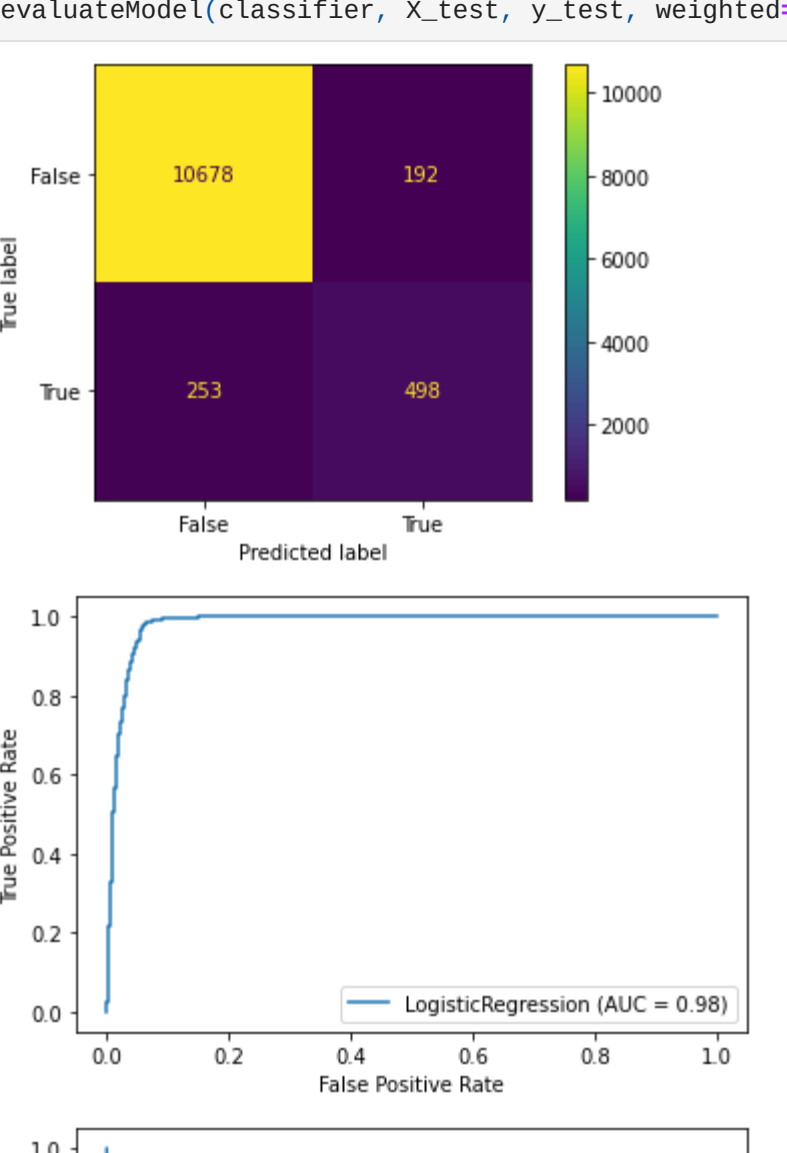
#Plotting class membership
import seaborn as sns
pl = sns.countplot(x=target)
plt.title('class distribution')
plt.show()
```



The class distribution bar chart shows the count of samples for each class. Class 3 (cover type 3) has the highest count, around 500,000.

Now let's see the same metrics without weighted averaging:

```
In [41]: evaluateModel(classifier, X_test, y_test, weighted=False)
```



The confusion matrix shows 107287 True Positives, 1795 False Positives, 2775 False Negatives, and 4346 True Negatives. The ROC curve shows an AUC of 0.98, and the Precision-Recall curve shows an AP of 0.70. The classification report shows an unweighted accuracy of 0.96, precision of 0.96, recall of 0.96, and F1 score of 0.96.

This is obvious because classes are severely imbalanced.

Now if we take the weighted average scoring, we will see very different scenario:

## Question 2

Should you get more training data?

## Answer

In my understanding, getting more data in such an imbalanced scenario may not help the classifier. We could go for class balancing, e.g., oversampling or under-sampling. But I'm not gonna do that (that's fairly easy, e.g., using SMOTHE from the imblearn library), as real-life data may not be balanced, e.g., credit card fraud detection, or money laundering use cases, severe imbalance is expected. To prove my hypothesis, let's try training the LR model on full dataset.

```
In [60]: COVER_TYPE = 3
features, target = getdataForCoverType3(dataset, COVER_TYPE, False)
X_train, X_test, y_train, y_test = train_test_split(features, target, test_size=0.2, random_state=42)

In [61]: import numpy as np

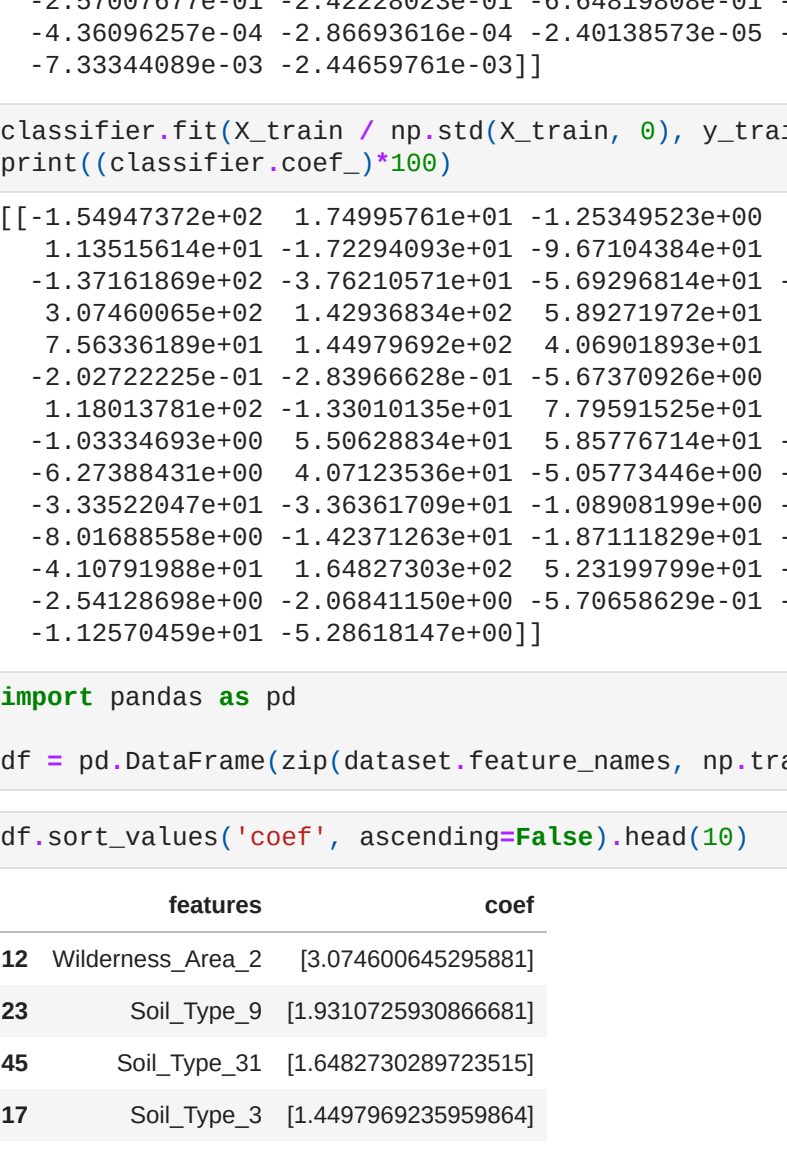
unique, counts = np.unique(target, return_counts=True)
dict(zip(unique, counts))

Out[61]: False: 54615, True: 3486)

In [62]: classifier.fit(X_train, y_train)

Out[62]: LogisticRegression(solver='liblinear')
```

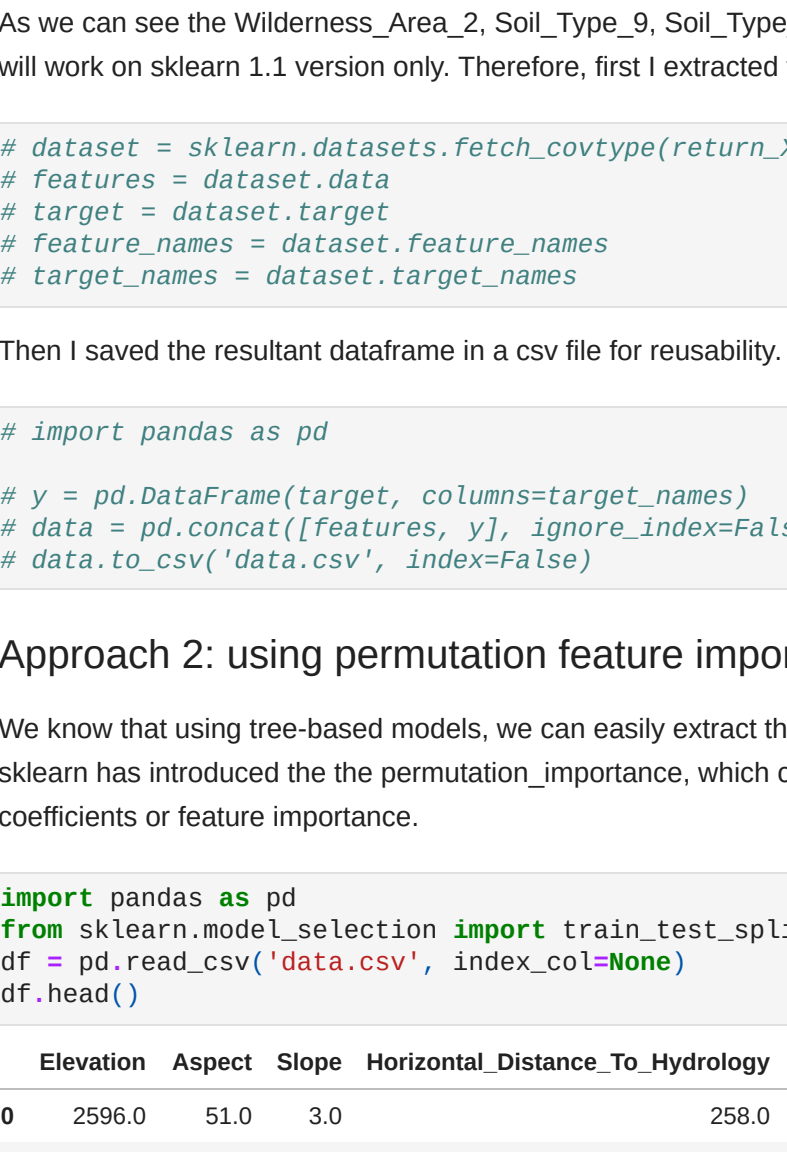
```
In [63]: evaluateModel(classifier, X_test, y_test, weighted=False)
```



The confusion matrix shows 10078 True Positives, 112 False Positives, 253 False Negatives, and 498 True Negatives. The ROC curve shows an AUC of 0.98, and the Precision-Recall curve shows an AP of 0.72. The classification report shows a weighted accuracy of 0.96, precision of 0.96, recall of 0.96, and F1 score of 0.96.

Now if we take the weighted average scoring, we will see very different scenario:

```
In [64]: evaluateModel(classifier, X_test, y_test, weighted=True)
```



The confusion matrix shows 10078 True Positives, 112 False Positives, 253 False Negatives, and 498 True Negatives. The ROC curve shows an AUC of 0.98, and the Precision-Recall curve shows an AP of 0.72. The classification report shows a weighted accuracy of 0.96, precision of 0.96, recall of 0.96, and F1 score of 0.96.

Overall, even training with all the available data did not help model perform well due to high data imbalance.

## Question 3

How would you decide which features to include in the deployed model?

```
In [29]: ### Write your code and explanation here.
```

## Answer

As not all the features could be of importance or relevant in the predictive modeling, as unimportant features will not only introduce intrinsic noises but also increase the computational complexity of the model. Therefore, it makes sense to identify the most impactful features based on their importance, followed by dropping less important features.

### Approach 1: w.r.t coefficients

The very basic thing I'll do is computing and examining the coefficients of the model on standardized parameters, the estimated coefficients, which will be around 1 across features. This then can be used as the feature importance.

```
In [30]: # The estimated coefficients will all be around 1:
print(classifier.coef_)

[[-7.78969185e-03 1.36235296e-03 1.07658099e-01 2.29777366e-03
 2.44569936e-03 -1.69712259e-04 7.35975987e-02 -3.25679706e-02
 5.02777376e-02 -2.96889625e-04 -1.21215252e+00 -1.90535866e-02
 9.97517882e-01 2.94581580e-01 -2.54903232e-01 1.44122848e+00
-2.69586376e-01 1.95502196e+00 -2.02871037e-01 8.89516511e-01
 5.03769698e-05 -1.12925389e-04 -6.18693992e-02 4.52315625e-01
1.51686484e-01 -2.51792650e-01 -6.19950246e-01 -2.17254850e-01
-1.92950886e-03 -7.78265559e-02 -2.10127526e-01 -5.68384667e-02
-6.07847471e-03 -1.20183575e-01 -1.96222846e-03 5.66811567e-02
-1.91784871e-01 -1.43508877e-01 -4.97215695e-04 -5.61565129e-02
-8.27231046e-03 -5.58283010e-02 -5.63976446e-01 -3.08464106e-01
-2.57707577e-01 -2.42228020e-01 -6.64819896e-01 -6.35915233e-02
-4.36896257e-04 -2.86693610e-04 -2.40138573e-05 -5.55165514e-03
-7.33344089e-03 -2.44659741e-03]]

In [18]: classifier.fit(X_train / np.std(X_train, 0), y_train)
print(classifier.coef_ * 100)

[[-1.54947372e+02 1.74995761e+01 -1.25349523e+00 4.17641255e+01
-1.13515614e+01 -1.72294093e+01 -9.07194384e+01 8.30884590e+01
-1.37161869e+02 -3.76210571e+01 -5.69296814e+01 -9.60462126e+00
3.07460605e+02 1.42936834e+02 5.89271972e+01 1.12366328e+02
-5.66336189e+01 -4.44077602e+02 4.06980189e+01 9.80797812e+01
-2.02722250e-01 -2.83966628e-01 -5.67370922e+00 1.93187259e+02
1.18913781e+02 -1.33010135e+01 7.79591525e+01 2.18578330e+01
-0.03334692e+00 5.50628034e+01 5.85776714e+01 -4.90902322e+00
-6.27388431e+00 4.07123536e+01 7.79591525e+00 -1.76797358e+01
-3.35522047e+01 -3.36361709e+01 -1.08986190e+00 -1.46283247e+01
-0.01680505e+00 -1.42371263e+01 -1.07111020e+01 1.52164001e+01
-4.10791988e+01 1.64887383e+02 5.23197099e+01 -1.39327246e+01
-2.54126698e+00 2.68241150e+00 -5.70658220e-01 -8.84928022e+00
-1.12570459e+01 -5.28618147e+00]]
```

```
In [19]: import pandas as pd

df = pd.DataFrame(zip(dataset.feature_names, np.transpose(classifier.coef_)), columns=['features', 'coef'])

In [20]: df.sort_values('coef', ascending=False).head(10)
```

	features	coef
23	Wilderness_Area_2	[3.0746006452956881]
12	Soil_Type_9	[1.0310725930666081]
45	Soil_Type_31	[1.6482730289723515]
17	Soil_Type_33	[1.4487969226969064]
13	Wilderness_Area_3	[1.4283683479409564]
24	Soil_Type_10	[1.1801378143985837]
15	Soil_Type_1	[1.1230632773649667]
19	Soil_Type_5	[0.8807978119024735]
7	Hillshade_Noon	[0.830874580762364]
26	Soil_Type_12	[0.7795915480171168]

As we can see the Wilderness\_Area\_2, Soil\_Type\_9, Soil\_Type\_31, Soil\_Type\_33, Wilderness\_Area\_3 are top-k features. However, this will work on sklearn 1.1 version only. Therefore, I'll extract feature and target names from the fetch\_covtype() method as follows:

```
In [ ]: dataset = sklearn.datasets.fetch_covtype(return_X_y=False, as_frame=True)
# features = dataset.data
# target = dataset.target
# feature_names = dataset.feature_names
# target_names = dataset.target_names

Then I saved the resultant dataframe in a csv file for reusability.

In [12]: # import pandas as pd

# y = pd.DataFrame(target, columns=target_names)
# data = pd.concat([features, y], ignore_index=False, axis=1).reset_index(drop=True)
# data.to_csv('data.csv', index=False)
```

### Approach 2: using permutation feature importance

We know that using tree-based models, we can easily extract the permutation feature importance (PFI). However, the recent version of sklearn introduced the permutation importance, which can be extracted via inspection wrapper for any model that provided coefficients or feature importance.

```
In [3]: import pandas as pd
from sklearn.model_selection import train_test_split, cross_val_score
df = pd.read_csv('data.csv', index_col=None)
df.head()
```

	Elevation	Aspect	Slope	Horizontal_Distance_To_Hydrology	Vertical_Distance_To_Hydrology	Horizontal_Distance_To_Roadways	Hillshade
0	2596.0	51.0	3.0		258.0	0.0	510.0
1	2590.0	56.0	2.0		212.0	-6.0	390.0
2	2804.0	139.0	9.0		268.0	65.0	3180.0
3	2785.0	155.0	18.0		143.0	-11.0	399.0
4	2595.0	45.0	2.0		252.0		3010.0

5 rows x 55 columns

```
In [4]: # Create a Pandas dataframe with all the Features
x = df.loc[:, df.columns != 'Cover_Type']
y = df['Cover_Type']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

In [ ]: from sklearn.inspection import permutation_importance
classifier = sklearn.linear_model.LogisticRegression(solver='liblinear')
classifier.fit(X_train, y_train)

pfi = permutation_importance(classifier, X_test, y_test, n_repeats=5, random_state=42)

In [27]: sorted_idx = pfi.importances.mean.argsort()
lr_importances = pd.Series(pfi.importances.mean, index=X_train.columns)

y_ticks = np.arange(0, len(X_train.columns))
fig, ax = plt.subplots(figsize=(15, 15))
ax.barh(y_ticks, lr_importances[sorted_idx])
ax.set_yticklabels(X_train.columns[sorted_idx])
ax.set_title('Feature Importances')
fig.tight_layout()
plt.show()
```



The Feature Importances bar chart shows the mean permutation importance for each feature. The top features are Hillshade\_3pm, Hillshade\_9am, Hillshade\_Noon, and Slope, all with importance values around 12.

This time, Hillshade\_3pm, Hillshade\_9am, Hillshade\_Noon, and Slope turn to be top-k features across different classes. This is very different compared to the first approach.

```
In [ ]: # from sklearn.inspection import DecisionBoundaryDisplay
# disp = DecisionBoundaryDisplay.from_estimator(classifier, X_test,
# response_method="predict", xlabel="Hillshade_3pm",
# ylabel="Elevation", alpha=0.5)
# disp.ax_.scatter(X_test[:, 0], X_test[:, 1], c=y_test, edgecolor="k")
# plt.show()

Approach 3: using SHAP

We already generated PFI using LR model, which provides a view on global feature importance and signifies how important a feature is for the model. However, PFI does not necessarily reflect the intrinsic predictive value of a feature by itself, since a feature that seems to have lower importance for an underfitted model could be important for a better-fitted model and the order in which a model observes features can affect the predictions. Therefore, I decided to employ SHAP to provide an explanation of the predictions of the model. SHAP uses Shapley values (SVs) as the measure of feature attribution and are based on coalition game theory. Further, as the order in which features are observed by the model matters, SHAP values explain the output of a function as a sum of the effects of each feature being introduced into a conditional expectation. If a feature has no or almost zero effect on the predicted value, it is expected to produce a SV of 0. If two features contribute equally to the prediction, SVs should have the same.
```

```
In [2]: import shap
# Initialize JavaScript visualization - use Jupyter notebook to see the interactive features of the plot
shap.initjs()
```

```
In [31]: # Create a TreeExplainer and extract shap values from it - will be used for plotting later
explainer = shap.Explainer(classifier, X_train, feature_names=X_train.columns)
shap_values = explainer.shap_values(X_test)

In [52]: shap_values = explainer(X_test)

In [32]: shap.summary_plot(shap_values, X_test, plot_type="bar")
```



The SHAP summary plot bar chart shows the mean(|SHAP value|) (average impact on model output) for each feature across different classes. The top features are Hillshade\_3pm, Hillshade\_9am, Hillshade\_Noon, and Slope, all with mean(|SHAP value|) around 12.

As seen, Hillshade\_3pm, Hillshade\_9am, Hillshade\_Noon, and Slope turn to be top-k features across different classes. This is somewhat identical to what we observed using the inspection wrapper.

### Approach 4: advanced techniques with AutoML

I used PyCaret - one of my favorite AutoML library. First, I'll compare performance between baseline classifiers and choose a top model whose performance will be evaluated w.r.t different performance metrics.

```
In [62]: from pycaret.classification import *
setup(df, target = 'Cover_Type')
```



0		Description	Value	
1		Target	Cover_Type	
2		Target Type	Multiclass	
3		Label Encoded	False	
4		Original Data	(581012, 55)	
5		Missing Values	False	
6		Numeric Features	44	
7		Categorical Features	10	
8		Ordinal Features	False	
9		High Cardinality Features	False	
10		High Cardinality Method	None	
11		Transformed Train Set	(406708, 54)	
12		Transformed Test Set	(174304, 54)	
13		Shuffle Train-Test	True	
14		Stratify Train-Test	False	
15		Fold Generator	StratifiedKFold	
16		Fold Number	10	
17		CPU Jobs	-1	
18		Use GPU	False	
19		Log Experiment	False	
20		Experiment Name	clf-default-name	
21		USI	6e5d	
22		Imputation Type	simple	
23		Iterative Imputation Numeric Model	None	
24		Numeric Imputer	mean	
25		Iterative Imputation Numeric Model	None	
26		Categorical Imputer	constant	
27		Iterative Imputation Categorical Model	None	
28		Unknown Categoricals Handling	least_frequent	
29		Normalize	False	
30		Normalize Method	None	
31		Transformation	False	
32		Transformation Method	None	
33		PCA	False	
34		PCA Method	None	
35		PCA Components	None	
36		Ignore Low Variance	False	
37		Combine Rare Levels	False	
38		Rare Level Threshold	None	
39		Numeric Binning	None	
40		Remove Outliers	False	
41		Outliers Threshold	None	
42		Remove Multicollinearity	False	
43		Multicollinearity Threshold	True	
44		Remove Perfect Collinearity	True	
45		Clustering	False	
46		Clustering Iteration	None	
47		Polynomial Features	False	
48		Polynomial Degree	None	
49		Trigonometry Features	False	
50		Polynomial Threshold	None	
51		Group Features	False	
52		Feature Selection Method	classic	
53		Feature Selection Threshold	None	
54		Feature Interaction	False	
55		Feature Ratio	False	
56		Interaction Threshold	None	
57		Fix Imbalance	False	
58		Fix Imbalance Method	SMOTE	
59		Fix Imbalance Method	SMOTE	
Out[62]:				
[<Pandas.io.formats.style.Styler at 0x7f80cc7127f8>],				
0		Elevation	Aspect Slope Horizontal_Distance_To_Hydrology \	
1	192545	2724	80.0 4.0 1.0 510.0	
2	335432	3182.0	285.0 13.0 1.0 42.0	
3	128469	2612.0	66.0 8.0 1.0 30.0	
4	384327	2933.0	84.0 4.0 1.0 693.0	
5	284912	3229.0	333.0 16.0 1.0 108.0	
...	...	...	...	
...	273280	2217.0	333.0 29.0 1.0 874.0	
...	46317	3234.0	275.0 11.0 1.0 525.0	
...	368693	3226.0	315.0 1.0 1.0 87.0	
...	242362	2837.0	291.0 3.0 1.0 124.0	
...	17454	2665.0	189.0 5.0 1.0 108.0	
...	...	...	...	
...	103958	16.0	0.0 0.0 1.0 1234.0	
...	335432	3.0	43.0 141.0 1.0 4462.0	
...	128469	182.0	240.0 196.0 1.0 212.0	
...	384327	229.0	238.0 128.0 1.0 2850.0	
...	284912	182.0	22.0 175.0 1.0 3141.0	
...	...	...	...	
...	273280	43.0	43.0 1434.0 1.0 1434.0	
...	46317	59.0	59.0 4142.0 1.0 4142.0	
...	368693	87.0	87.0 153.0 1.0 693.0	
...	242362	22.0	22.0 1224.0 1.0 1224.0	
...	17454	27.0	27.0 1590.0 1.0 1590.0	
...	...	...	...	
...	Hillshade_9am	Hillshade Noon Hillshade_3pm \		
0	103958	226.0	233.0 141.0 0.0 0.0	
1	335432	183.0	240.0 196.0 1.0 182.0	
2	128469	229.0	224.0 128.0 1.0 525.0	
3	384327	226.0	232.0 141.0 1.0 2850.0	
4	284912	182.0	218.0 175.0 1.0 3141.0	
...	...	...	...	
...	273280	141.0	188.0 179.0 1.0 874.0	
...	46317	189.0	243.0 193.0 1.0 124.0	
...	368693	216.0	237.0 159.0 1.0 87.0	
...	242362	110.0	213.0 233.0 1.0 1224.0	
...	17454	220.0	243.0 159.0 1.0 1590.0	
...	...	...	...	
...	Horizontal_Distance_To_Fire_Points ... Soil_Type_30_0_0 \			
...	103958	2329.0	1.0 1.0 1.0 182.0	
...	335432	1782.0	...	1.0
...	128469	430.0	...	1.0
...	384327	2791.0	...	1.0
...	284912	4195.0	...	1.0
...	...	...	...	...
...	273280	698.0	...	1.0
...	46317	2481.0	...	1.0
...	368693	1746.0	...	0.0
...	242362	1289.0	...	1.0
...	17454	6584.0	...	1.0
...	...	...	...	...
...	Soil_Type_31_1_0 Soil_Type_32_0_0 Soil_Type_33_0_0 \			
...	103958	0.0	1.0 1.0 1.0	
...	335432	0.0	1.0 1.0 1.0	
...	128469	0.0	1.0 1.0 1.0	
...	384327	0.0	1.0 1.0 1.0	
...	284912	0.0	1.0 1.0 1.0	
...	...	...	...	
...	273280	0.0	1.0 1.0 1.0	
...	46317	0.0	1.0 1.0 1.0	
...	368693	0.0	1.0 1.0 1.0	
...	242362	0.0	1.0 1.0 1.0	
...	17454	0.0	1.0 1.0 1.0	
...	...	...	...	
...	Soil_Type_34_0_0 Soil_Type_35_1_0 Soil_Type_36_0_0 \			
...	103958	1.0	0.0 1.0 1.0	
...	335432	0.0	1.0 1.0 1.0	
...	128469	1.0	0.0 1.0 1.0	
...	384327	1.0	0.0 1.0 1.0	
...	284912	1.0	0.0 1.0 1.0	
...	...	...	...	
...	273280	1.0	0.0 1.0 1.0	
...	46317	0.0	1.0 1.0 1.0	
...	368693	1.0	0.0 1.0 1.0	
...	242362	1.0	0.0 1.0 1.0	
...	17454	1.0	0.0 1.0 1.0	
...	...	...	...	
...	Soil_Type_37_1_0 Soil_Type_38_1_0 Soil_Type_39_0_0 \			
...	103958	0.0	0.0 0.0 1.0	
...	335432	0.0	0.0 1.0 1.0	
...	128469	0.0	0.0 1.0 1.0	
...	384327	0.0	0.0 1.0 1.0	
...	284912	0.0	0.0 1.0 1.0	
...	...	...	...	
...	273280	0.0	0.0 1.0 1.0	
...	46317	0.0	0.0 1.0 1.0	
...	368693	0.0	0.0 1.0 1.0	
...	242362	0.0	0.0 1.0 1.0	
...	17454	0.0	0.0 1.0 1.0	
...	...	...	...	
[174304 rows x 54 columns],				
0		Elevation	Aspect Slope Horizontal_Distance_To_Hydrology \	
1	2596.0	51.0	3.0 1.0 258.0	
2	2884.0	139.0	9.0 1.0 268.0	
3	2785.0	155.0	18.0 1.0 242.0	
4	2595.0	45.0	2.0 1.0 693.0	
...	...	...	...	
...	581007	2396.0	153.0 20.0 1.0 85.0	
...	581091	2391.0	152.0 19.0 1.0 124.0	
...	581099	2386.0	159.0 17.0 1.0 69.0	
...	581010	2384.0	170.0 15.0 1.0 60.0	
...	581011	2383.0	105.0 13.0 1.0 60.0	
...	...	...	...	
...	Vertical_Distance_To_Hydrology Horizontal_Distance_To_Roadways \			
0	103958	16.0	0.0 0.0 1.0 510.0	
1	335432	3.0	43.0 141.0 1.0 4462.0	
2	128469	182.0	240.0 196.0 1.0 212.0	
3	384327	229.0	238.0 128.0 1.0 2850.0	
4	284912	182.0	22.0 175.0 1.0 3141.0	
...	...	...	...	
...	273280	43.0	43.0 1434.0 1.0 1434.0	
...	46317	59.0	59.0 4142.0 1.0 4142.0	
...	368693	87.0	87.0 153.0 1.0 693.0	
...	242362	22.0	22.0 1224.0 1.0 1224.0	
...	17454	27.0	27.0 1590.0 1.0 1590.0	
...	...	...	...	
...	Hillshade_9am Hillshade Noon Hillshade_3pm \			
0	103958	226.0	233.0 141.0 0.0 0.0	
1	335432	183.0	240.0 196.0 1.0 182.0	
2	128469	229.0	224.0 128.0 1.0 525.0	
3	384327	226.0	232.0 141.0 1.0 2850.0	
4	284912	182.0	218.0 175.0 1.0 3141.0	
...	...	...	...	
...	273280	141.0	188.0 179.0 1.0 874.0	
...	46317	189.0	243.0 193.0 1.0 124.0	
...	368693	216.0	237.0 159.0 1.0 87.0	
...	242362	110.0	213.0 233.0 1.0 1224.0	
...	17454	220.0	243.0 159.0 1.0 1590.0	
...	...	...	...	
...	Horizontal_Distance_To_Fire_Points ... Soil_Type_30_0_0 \			
...	103958	2329.0	1.0 1.0 1.0 182.0	
...	335432	1782.0	...	1.0
...	128469	430.0	...	1.0
...	384327	2791.0	...	1.0
...	284912	4195.0	...	1.0
...	...	...	...	...
...	273280	698.0	...	1.0
...	46317	2481.0	...	1.0
...	368693	1746.0	...	0.0
...	242362	1289.0	...	1.0
...	17454	6584.0	...	1.0
...	...	...	...	...
...	Soil_Type_31_1_0 Soil_Type_32_0_0 Soil_Type_33_0_0 \			
...	103958	0.0	1.0 1.0 1.0	
...	335432	0.0	1.0 1.0 1.0	
...	128469	0.0	1.0 1.0 1.0	
...	384327	0.0	1.0 1.0 1.0	
...	284912	0.0	1.0 1.0 1.0	
...	...	...	...	
...	273280	0.0	1.0 1.0 1.0	
...	46317	0.0	1.0 1.0 1.0	
...	368693	0.0	1.0 1.0 1.0	
...	242362	0.0	1.0 1.0 1.0	
...	17454	0.0	1.0 1.0 1.0	
...	...	...	...	
...	Soil_Type_34_0_0 Soil_Type_35_1_0 Soil_Type_36_0_0 \			
...	103958	1.0	0.0 1.0 1.0	
...	335432	0.0	1.0 1.0 1.0	
...	128469	1.0	0.0 1.0 1.0	
...	384327	1.0	0.0 1.0 1.0	
...	284912	1.0	0.0 1.0 1.0	
...	...	...	...	
...	273280	1.0	0.0 1.0 1.0	
...	46317	0.0	1.0 1.0 1.0	
...	368693	1.0	0.0 1.0 1.0	
...	242362	1.0	0.0 1.0 1.0	
...	17454	1.0	0.0 1.0 1.0	
...	...	...	...	
...	Soil_Type_37_1_0 Soil_Type_38_1_0 Soil_Type_39_0_0 \			
...	103958	0.0	0.0 0.0 1.0	
...	335432	0.0	0.0 1.0 1.0	
...	128469	0.0	0.0 1.0 1.0	
...	384327	0.0	0.0 1.0 1.0	
...	284912	0.0	0.0 1.0 1.0	
...	...	...	...	
...	273280	0.0	0.0 1.0 1.0	
...	46317	0.0	0.0 1.0 1.0	
...	368693	0.0	0.0 1.0 1.0	
...	242362	0.0	0.0 1.0 1.0	
...	17454	0.0	0.0 1.0 1.0	
...	...	...	...	
[174304 rows x 54 columns],				
0		Elevation	Aspect Slope Horizontal_Distance_To_Hydrology \	
1	2596.0	51.0	3.0 1.0 258.0	
2	2884.0	139.0	9.0 1.0 268.0	
3	2785.0	155.0	18.0 1.0 242.0	
4	2595.0	45.0	2.0 1.0 693.0	
...	...	...	...	
...	581007	2396.0	153.0 20.0 1.0 85.0	
...	581091	2391.0	152.0 19.0 1.0 124.0	
...	581099	2386.0	159.0 17.0 1.0 69.0	
...	581010	2384.0	170.0 15.0 1.0 60.0	
...	581011	2383.0	105.0 13.0 1.0 60.0	
...	...	...	...	
...	Vertical_Distance_To_Hydrology Horizontal_Distance_To_Roadways \			
0	103958	16.0	0.0 0.0 1.0 510.0	
1	335432	3.0	43.0 141.0 1.0 4462.0	
2	128469	182.0	240.0 196.0 1.0 212.0	
3	384327	229.0	238.0 128.0 1.0 2850.0	
4	284912	182.0	22.0 175.0 1.0 3141.0	
...	...	...	...	
...	273280	43.0	43.0 1434.0 1.0 1434.0	
...	46317	59.0	59.0 4142.0 1.0 4142.0	
...	368693	87.0	87.0 153.0 1.0 693.0	
...	242362	22.0	22.0 1224.0 1.0 1224.0	
...	17454	27.0	27.0 1590.0 1.0 1590.0	
...	...	...	...	
...	Hillshade_9am Hillshade Noon Hillshade_3pm \			
0	103958	226.0	233.0 141.0 0.0 0.0	
1	335432	183.0	240.0 196.0 1.0 182.0	
2	128469	229.0	224.0 128.0 1.0 525.0	
3	384327	226.0	232.0 141.0 1.0 2850.0	
4	284912	182.0	218.0 175.0 1.0 3141.0	
...	...	...	...	
...	273280	141.0	188.0 179.0 1.0 874.0	
...	46317	189.0	243.0 193.0 1.0 124.0	
...	368693	216.0	237.0 159.0 1.0 87.0	
...	242362	110.0	213.0 233.0 1.0 1224.0	
...	17454	220.0	243.0 159.0 1.0 1590.0	
...	...	...	...	
...	Horizontal_Distance_To_Fire_Points ... Soil_Type_30_0_0 \			
...	103958	2329.0	1.0 1.0 1.0 182.0	
...	335432	1782.0	...	1.0
...	128469	430.0	...	1.0



```
581899      1.0      1.0      1.0
581910      1.0      1.0      1.0
581911      0.0      1.0      1.0

Soil_Type_35_1_0  Soil_Type_35_1_0  Soil_Type_36_0_0  \
0      1.0      0.0      1.0
1      1.0      0.0      1.0
2      1.0      0.0      1.0
3      1.0      0.0      1.0
4      1.0      0.0      1.0
...      ...      ...      ...
581897      1.0      0.0      1.0
581898      1.0      0.0      1.0
581899      1.0      0.0      1.0
581910      1.0      0.0      1.0
581911      1.0      0.0      1.0

Soil_Type_37_1_0  Soil_Type_38_1_0  Soil_Type_39_0_0
0      0.0      0.0      1.0
1      0.0      0.0      1.0
2      0.0      0.0      1.0
3      0.0      0.0      1.0
4      0.0      0.0      1.0
...      ...      ...      ...
581897      0.0      0.0      1.0
581898      0.0      0.0      1.0
581899      0.0      0.0      1.0
581910      0.0      0.0      1.0
581911      0.0      0.0      1.0

[581812 rows x 54 columns],
0.0,
None,
'lightgbm',
False,
10,
None,
57439      2
171107      2
411248      1
225267      1
67415      2
287843      -
462415      2
349350      3
429600      1
109060      1
Name: Cover_Type, Length: 466768, dtype: int64,
'lightgbm'
```

```
In [63]: # This model will be used to compare all the model along with the cross validation
best = compare_models()
```

	Model	Accuracy	AUC	Recall	Prec.	F1	Kappa	MCC	TT (Sec)
knn	K Neighbors Classifier	0.9635	0.9959	0.9202	0.9635	0.9635	0.9414	0.9414	1.3810
rf	Random Forest Classifier	0.9503	0.9944	0.8974	0.9504	0.9500	0.9198	0.9200	2.4210
et	Extra Trees Classifier	0.9492	0.9938	0.9019	0.9493	0.9499	0.9191	0.9183	2.9000
dt	Decision Tree Classifier	0.9317	0.9450	0.8909	0.9318	0.9318	0.8905	0.8905	0.8010
lightgbm	Light Gradient Boosting Machine	0.8516	0.9523	0.8091	0.8527	0.8515	0.7606	0.7610	2.1520
gbc	Gradient Boosting Classifier	0.7725	0.9085	0.6495	0.7718	0.7687	0.6282	0.6289	73.2890
ridge	Ridge Classifier	0.7007	0.0000	0.3716	0.6787	0.6779	0.5032	0.5085	0.2610
lr	Logistic Regression	0.6840	0.6353	0.3465	0.6503	0.6572	0.4648	0.4680	17.8680
lda	Linear Discriminant Analysis	0.6794	0.8369	0.5763	0.6941	0.6825	0.4963	0.4977	0.8170
svm	SVM - Linear Kernel	0.5326	0.0000	0.3238	0.6019	0.4838	0.2672	0.3134	19.4620
ada	Ada Boost Classifier	0.5190	0.6042	0.3673	0.6222	0.5348	0.2993	0.3150	2.4090
dummy	Dummy Classifier	0.4875	0.5000	0.1429	0.2377	0.3195	0.0000	0.0000	0.2410
nb	Naive Bayes	0.4571	0.7814	0.5877	0.6451	0.4153	0.2631	0.3076	0.2890
qda	Quadratic Discriminant Analysis	0.4542	0.5375	0.3036	0.4383	0.3569	0.0761	0.1080	0.3470

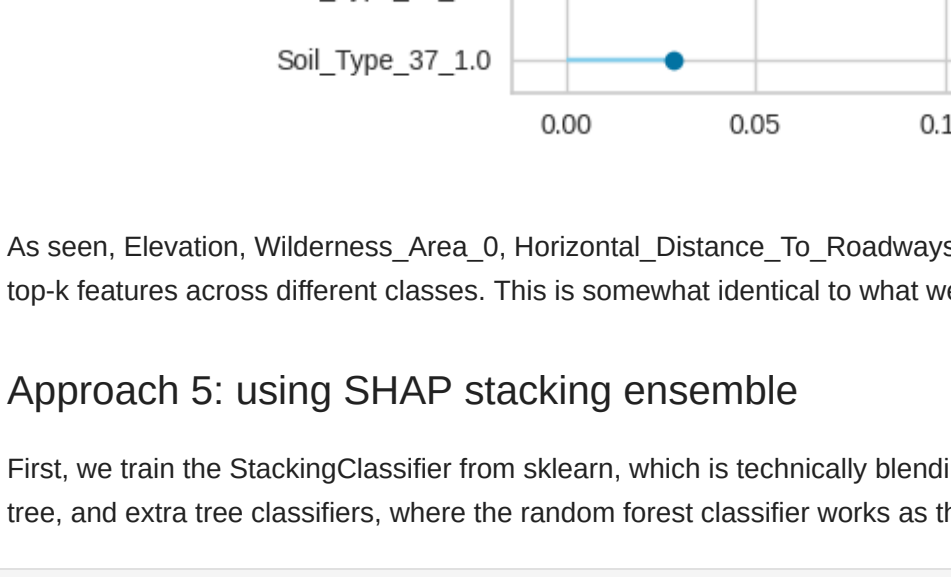
```
In [71]: # Build random forest classifier.
# I picked it even though it tooked 2 due to the fact that we can extract feature importance
rf=create_model('rf')
```

	Accuracy	AUC	Recall	Prec.	F1	Kappa	MCC
Fold							
0	0.9493	0.9942	0.8912	0.9493	0.9489	0.9182	0.9183
1	0.9509	0.9945	0.8992	0.9510	0.9506	0.9207	0.9210
2	0.9502	0.9941	0.8978	0.9502	0.9498	0.9196	0.9198
3	0.9492	0.9944	0.8899	0.9493	0.9488	0.9180	0.9182
4	0.9512	0.9947	0.9045	0.9514	0.9510	0.9213	0.9216
5	0.9507	0.9944	0.8973	0.9509	0.9504	0.9205	0.9207
6	0.9502	0.9945	0.8990	0.9503	0.9500	0.9197	0.9199
7	0.9511	0.9947	0.9088	0.9511	0.9508	0.9211	0.9213
8	0.9491	0.9939	0.8829	0.9492	0.9488	0.9179	0.9181
9	0.9510	0.9949	0.9038	0.9511	0.9508	0.9219	0.9212
Mean	0.9503	0.9944	0.8974	0.9504	0.9500	0.9198	0.9200
Std	0.0008	0.0003	0.0073	0.0008	0.0003	0.0013	0.0013

```
In [72]: # Whenever we compare different models or build a due to the fact that we can extract feature importance
# hyperparameter values. Hence, we need to tune our model to get better performance
tuned_rf=tune_model(rf)
```

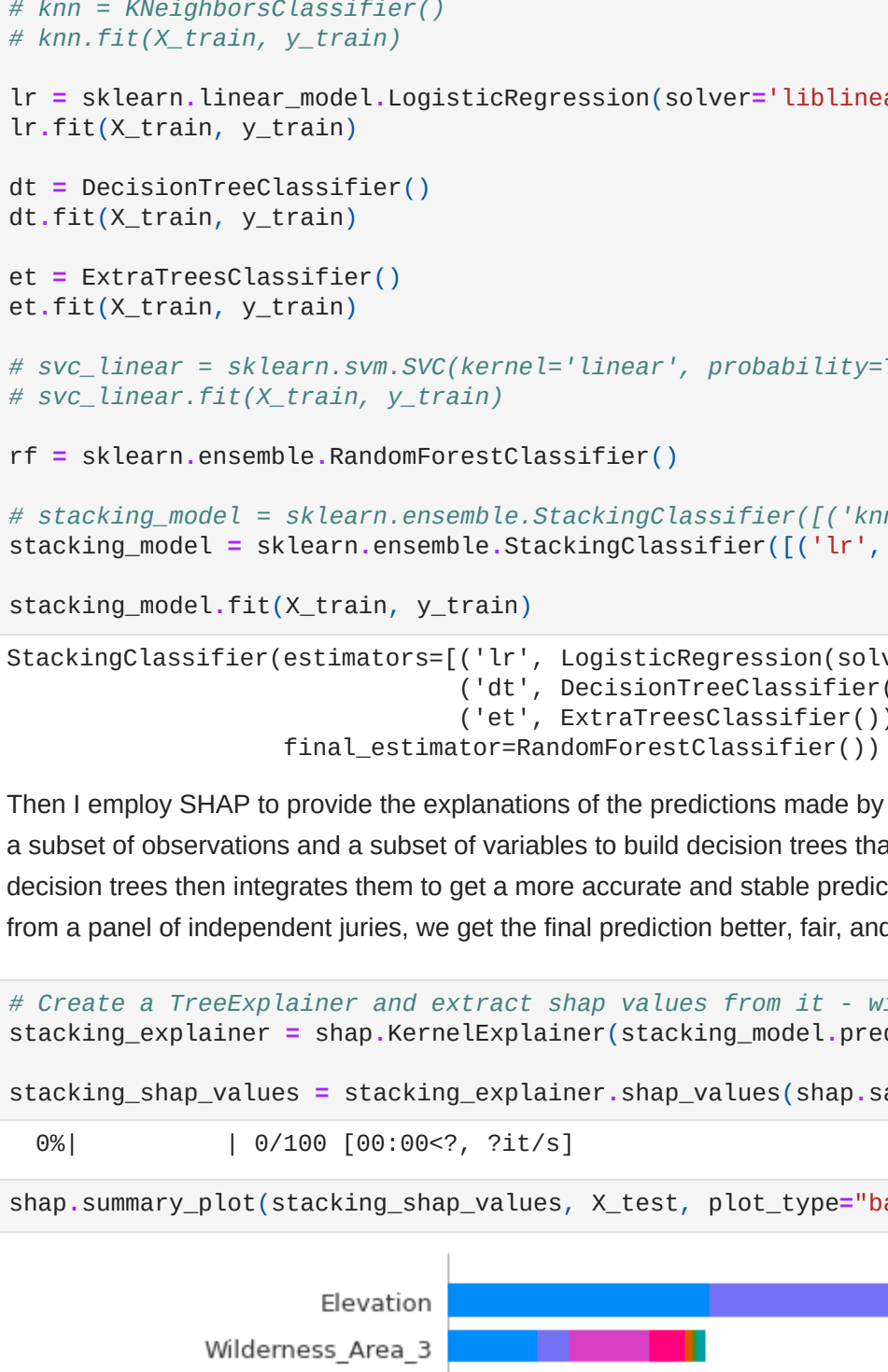
	Accuracy	AUC	Recall	Prec.	F1	Kappa	MCC
Fold							
0	0.5197	0.7917	0.7069	0.6875	0.5476	0.3512	0.3782
1	0.5218	0.7910	0.7056	0.6831	0.5497	0.3515	0.3772
2	0.5230	0.7889	0.7110	0.6827	0.5512	0.3526	0.3778
3	0.5188	0.7885	0.7043	0.6816	0.5469	0.3486	0.3745
4	0.5291	0.7894	0.7145	0.6900	0.5575	0.3610	0.3872
5	0.5212	0.7890	0.7062	0.6895	0.5502	0.3530	0.3799
6	0.5177	0.7889	0.7017	0.6850	0.5468	0.3490	0.3758
7	0.5241	0.7876	0.7104	0.6872	0.5519	0.3553	0.3818
8	0.5294	0.7879	0.7081	0.6826	0.5491	0.3510	0.3768
9	0.5183	0.7915	0.7106	0.6826	0.5461	0.3484	0.3748
Mean	0.5214	0.7884	0.7079	0.6852	0.5497	0.3521	0.3784
Std	0.0032	0.0014	0.0036	0.0030	0.0032	0.0036	0.0036

```
In [73]: plot_model(tuned_rf, plot='boundary')
```

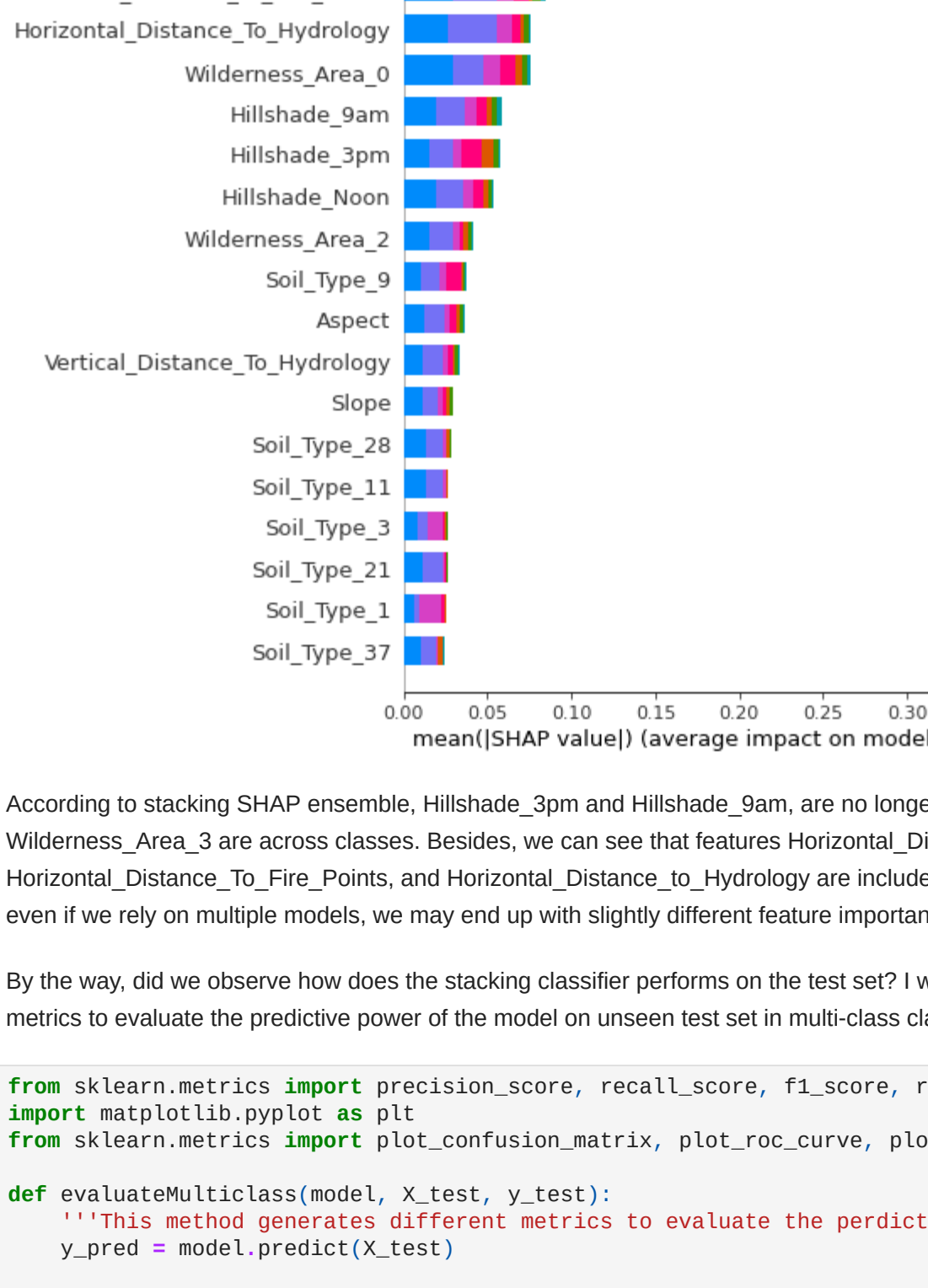


It seems the RF model is able to create a clear separation of samples w.r.t decision boundary.

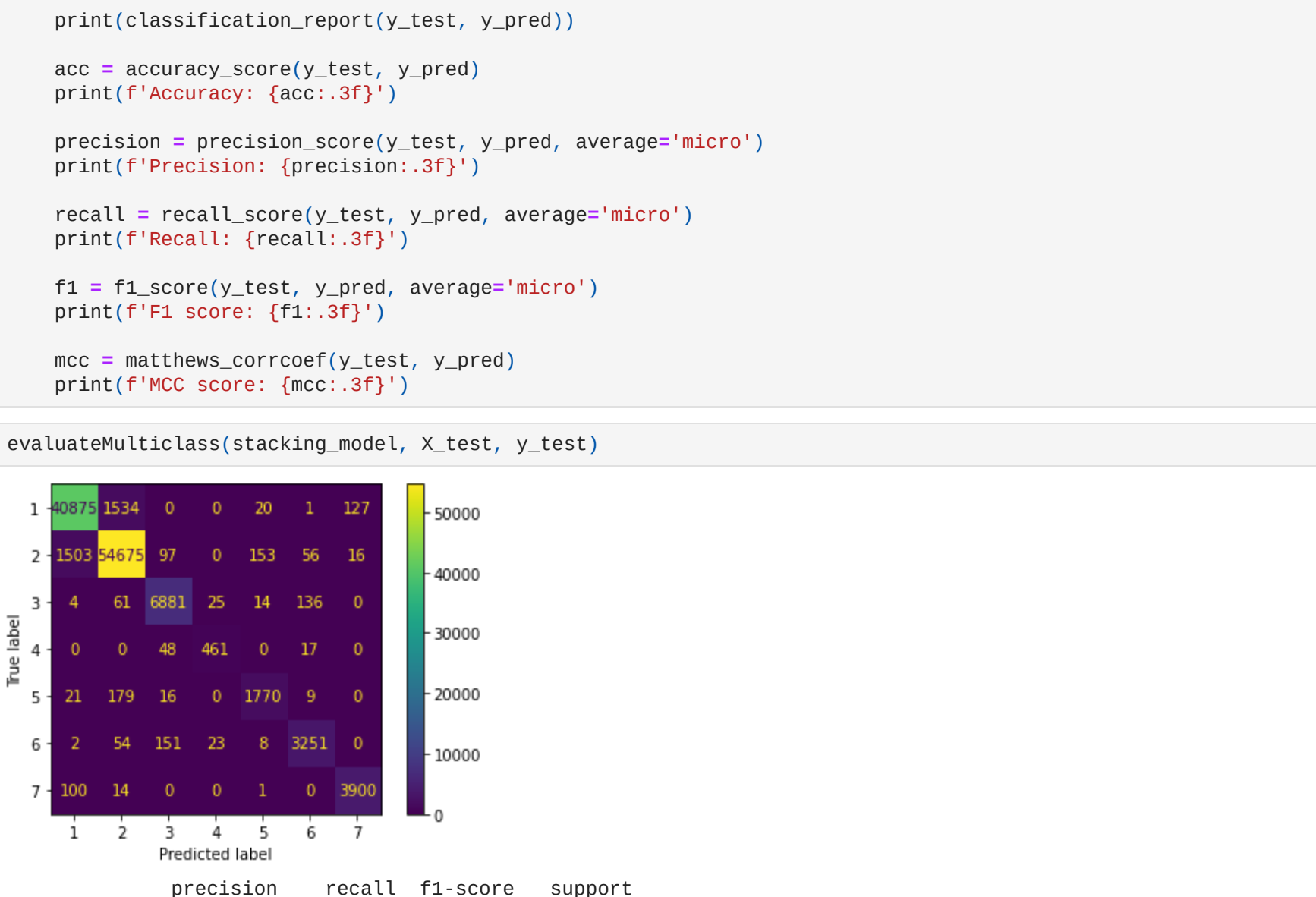
```
In [74]: plot_model(tuned_rf, plot='auc', scale = 1)
```



```
In [75]: plot_model(tuned_rf, plot='confusion_matrix', plot_kwards = {'percent' : True})
```



```
In [76]: plot_model(tuned_rf, plot='feature')
```



As seen, Elevation, Wilderness\_Area\_0, Horizontal\_Distance\_To\_Roadways, Wilderness\_Area\_0, and Wilderness\_Area\_2 turn to be are top-k features across different classes. Besides, we can see that features Horizontal\_Distance\_To\_Roadways, Horizontal\_Distance\_To\_Fire\_Points, and Horizontal\_Distance\_To\_Hydrology are included in top-5 important feature list. In summary, even if we rely on multiple models, we may end up with slightly different feature importance.

## Approach 5: using SHAP stacking ensemble

First, we train the StackingClassifier from sklearn, which is technically blending or stacking predictions from logistic regression, decision tree, and extra tree classifiers, where the random forest classifier works as the base estimator.

```
In [6]: import shap
import sklearn
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import ExtraTreesClassifier

# knn = KNeighborsClassifier()
# knn.fit(X_train, y_train)

lr = sklearn.linear_model.LogisticRegression(solver='liblinear')
lr.fit(X_train, y_train)

dt = DecisionTreeClassifier()
dt.fit(X_train, y_train)

et = ExtraTreesClassifier()
et.fit(X_train, y_train)

# svc_linear = sklearn.svm.SVC(kernel='linear', probability=True)
# svc_linear.fit(X_train, y_train)

rf = sklearn.ensemble.RandomForestClassifier()

# stacking_model = sklearn.ensemble.StackingClassifier([('lr', knn), ('svc', svc_linear)], rf)
stacking_model = sklearn.ensemble.StackingClassifier([('lr', lr), ('dt', dt), ('et', et)], rf)
stacking_model.fit(X_train, y_train)

Out[6]: StackingClassifier(estimators=[('lr', LogisticRegression(solver='liblinear')),
                                     ('dt', DecisionTreeClassifier()),
                                     ('et', ExtraTreesClassifier())],
                           final_estimator=RandomForestClassifier())
```

Then I employ SHAP to find the explanations of the predictions made by StackingClassifier. RF is an ensemble technique that takes a subset of observations and a subset of variables to build decision trees that creates an ensemble of decision trees. RF first builds decision trees then integrates them to get a more accurate and stable prediction. This is a direct consequence since by maximum voting from a panel of independent juries, we get the final prediction better, fair, and trustworthy than the best jury.

```
In [10]: # Create a TreeExplainer and extract shap values from it - will be used for plotting later
stacking_explainer = shap.KernelExplainer(stacking_model.predict_proba, shap.sample(X_test, 100))
stacking_shap_values = stacking_explainer.shap_values(shap.sample(X_test, 100))
0% | 0/100 [00:00<?, ?it/s]
```

```
In [11]: shap.summary_plot(stacking_shap_values, X_test, plot_type='bar')
```



According to stacking SHAP ensemble, Hillshade\_3pm and Hillshade\_9am, are no longer top-2 features, but Elevation and Wilderness\_Area\_3 are across classes. Besides, we can see that features Horizontal\_Distance\_To\_Roadways, Horizontal\_Distance\_To\_Fire\_Points, and Horizontal\_Distance\_To\_Hydrology are included in top-5 important feature list. In summary, even if we rely on multiple models, we may end up with slightly different feature importance.

By the way, did we observe how does the stacking classifier performs on the test set? I wrote a method, which generates different metrics to evaluate the predictive power of the model on unseen test set in multi-class classification setting.

```
In [24]: from sklearn.metrics import precision_score, recall_score, f1_score, roc_auc_score, accuracy_score, classification_report
import matplotlib.pyplot as plt
from sklearn.metrics import plot_confusion_matrix, plot_roc_curve, plot_precision_recall_curve

def evaluateMulticlass(model, X_test, y_test):
    '''This method generates different metrics to evaluate the predictive power of the model on unseen test set.
    y_pred = model.predict(X_test)

    plot_confusion_matrix(model, X_test, y_test)
    plt.show()

    print(classification_report(y_test, y_pred))

    acc = accuracy_score(y_test, y_pred)
    print(f'Accuracy: (acc:.3f)')

    precision = precision_score(y_test, y_pred, average='micro')
    print(f'Precision: (precision:.3f)')

    recall = recall_score(y_test, y_pred, average='micro')
    print(f'Recall: (recall:.3f)')

    f1 = f1_score(y_test, y_pred, average='micro')
    print(f'F1 score: (f1:.3f)')

    mcc = matthews_corcoef(y_test, y_pred)
    print(f'MCC score: (mcc:.3f)')
```

```
In [25]: evaluateMulticlass(stacking_model, X_test, y_test)
```



According to stacking SHAP ensemble, Hillshade\_3pm and Hillshade\_9am, are no longer top-2 features, but Elevation and Wilderness\_Area\_3 are across classes. Besides, we can see that features Horizontal\_Distance\_To\_Roadways, Horizontal\_Distance\_To\_Fire\_Points, and Horizontal\_Distance\_To\_Hydrology are included in top-5 important feature list. In summary, even if we rely on multiple models, we may end up with slightly different feature importance.

By the way, did we observe how does the stacking classifier performs on the test set? I wrote a method, which generates different metrics to evaluate the predictive power of the model on unseen test set in multi-class classification setting.

```
In [24]: from sklearn.metrics import precision_score, recall_score, f1_score, roc_auc_score, accuracy_score, classification_report
import matplotlib.pyplot as plt
from sklearn.metrics import plot_confusion_matrix, plot_roc_curve, plot_precision_recall_curve

def evaluateMulticlass(model, X_test, y_test):
    '''This method generates different metrics to evaluate the predictive power of the model on unseen test set.
    y_pred = model.predict(X_test)

    plot_confusion_matrix(model, X_test, y_test)
    plt.show()

    print(classification_report(y_test, y_pred))

    acc = accuracy_score(y_test, y_pred)
    print(f'Accuracy: (acc:.3f)')

    precision = precision_score(y_test, y_pred, average='micro')
    print(f'Precision: (precision:.3f)')

    recall = recall_score(y_test, y_pred, average='micro')
    print(f'Recall: (recall:.3f)')

    f1 = f1_score(y_test, y_pred, average='micro')
    print(f'F1 score: (f1:.3f)')

    mcc = matthews_corcoef(y_test, y_pred)
    print(f'MCC score: (mcc:.3f)')
```

```
In [25]: evaluateMulticlass(stacking_model, X_test, y_test)
```



According to stacking SHAP ensemble, Hillshade\_3pm and Hillshade\_9am, are no longer top-2 features, but Elevation and Wilderness\_Area\_3 are across classes. Besides, we can see that features Horizontal\_Distance\_To\_Roadways, Horizontal\_Distance\_To\_Fire\_Points, and Horizontal\_Distance\_To\_Hydrology are included in top-5 important feature list. In summary, even if we rely on multiple models, we may end up with slightly different feature importance.

By the way, did we observe how does the stacking classifier performs on the test set? I wrote a method, which generates different metrics to evaluate the predictive power of the model on unseen test set in multi-class classification setting.

```
In [24]: from sklearn.metrics import precision_score, recall_score, f1_score, roc_auc_score, accuracy_score, classification_report
import matplotlib.pyplot as plt
from sklearn.metrics import plot_confusion_matrix, plot_roc_curve, plot_precision_recall_curve

def evaluateMulticlass(model, X_test, y_test):
    '''This method generates different metrics to evaluate the predictive power of the model on unseen test set.
    y_pred = model.predict(X_test)

    plot_confusion_matrix(model, X_test, y_test)
    plt.show()

    print(classification_report(y_test, y_pred))

    acc = accuracy_score(y_test, y_pred)
    print(f'Accuracy: (acc:.3f)')

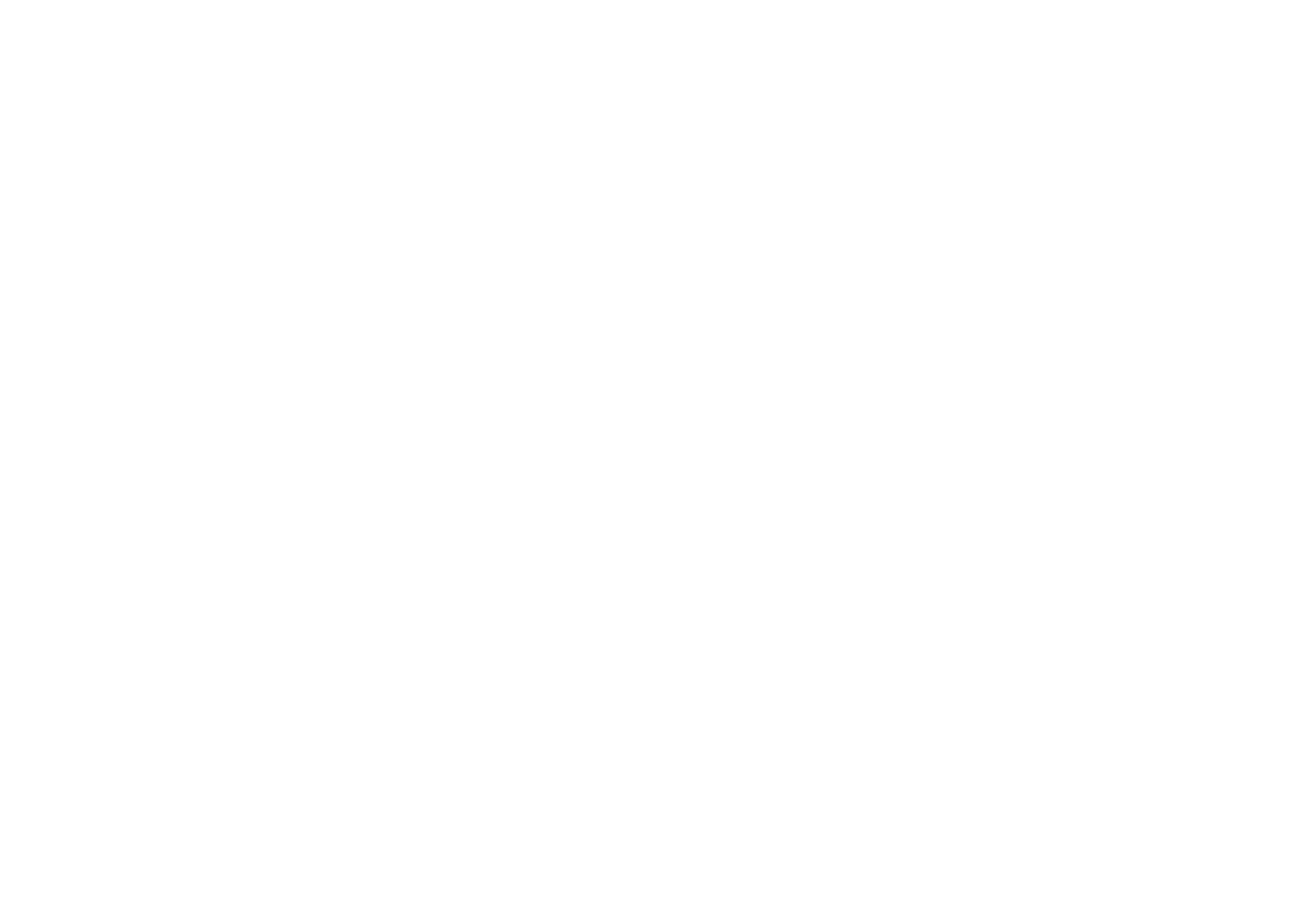
    precision = precision_score(y_test, y_pred, average='micro')
    print(f'Precision: (precision:.3f)')

    recall = recall_score(y_test, y_pred, average='micro')
    print(f'Recall: (recall:.3f)')

    f1 = f1_score(y_test, y_pred, average='micro')
    print(f'F1 score: (f1:.3f)')

    mcc = matthews_corcoef(y_test, y_pred)
    print(f'MCC score: (mcc:.3f)')
```

```
In [25]: evaluateMulticlass(stacking_model, X_test, y_test)
```



According to stacking SHAP ensemble, Hillshade\_3pm and Hillshade\_9am, are no longer top-2 features, but Elevation and Wilderness\_Area\_3 are across classes. Besides, we can see that features Horizontal\_Distance\_To\_Roadways, Horizontal\_Distance\_To\_Fire\_Points, and Horizontal\_Distance\_To\_Hydrology are included in top-5 important feature list. In summary, even if we rely on multiple models, we may end up with slightly different feature importance.

By the way, did we observe how does the stacking classifier performs on the test set? I wrote a method, which generates different metrics to evaluate the predictive power of the model on unseen test set in multi-class classification setting.

```
In [24]: from sklearn.metrics import precision_score, recall_score, f1_score, roc_auc_score, accuracy_score, classification_report
import matplotlib.pyplot as plt
from sklearn.metrics import plot_confusion_matrix, plot_roc_curve, plot_precision_recall_curve

def evaluateMulticlass(model, X_test, y_test):
    '''This method generates different metrics to evaluate the predictive power of the model on unseen test set.
    y_pred = model.predict(X_test)

    plot_confusion_matrix(model, X_test, y_test)
    plt.show()

    print(classification_report(y_test, y_pred))

    acc = accuracy_score(y_test, y_pred)
    print(f'Accuracy: (acc:.3f)')

    precision = precision_score(y_test, y_pred, average='micro')
    print(f'Precision: (precision:.3f)')

    recall = recall_score(y_test, y_pred, average='micro')
    print(f'Recall: (recall:.3f)')

    f1 = f1_score(y_test, y_pred, average='micro')
    print(f'F1 score: (f1:.3f)')

    mcc = matthews_corcoef(y_test, y_pred)
    print(f'MCC score: (mcc:.3f)')
```

```
In [25]: evaluateMulticlass(stacking_model, X_test, y_test)
```

