

11 Supervised learning with logistic regression

Nandi Christmas

```
In [49]: import numpy as np
import pandas as pd
from sklearn.preprocessing import StandardScaler
df_loan = pd.read_csv(
    filepath_or_url = "D:\\\\loan_data.csv",
    engine = "pyarrow"
)

Out[49]:
```

	loan_status	partition	person_education_2_moderate	person_education_3_high	person_home_ownership_3_high	loan_intent_2_moderate	loan_intent_3_high
0	1	test	0	0	1	1	0
1	0	test	1	0	0	0	0
2	1	train	1	0	0	0	1
3	1	train	0	1	1	1	0
4	1	test	0	0	0	1	0
...
44995	1	test	0	0	1	1	0
44996	1	validation	0	0	1	1	0
44997	1	train	0	0	1	0	1
44998	1	train	0	1	1	0	0
44999	1	train	1	0	1	0	1

45000 rows x 7 columns

```
In [50]: list_predictors = [
    'person_education_2_moderate',
    'person_education_3_high',
    'person_home_ownership_3_high',
    'loan_intent_2_moderate',
    'loan_intent_3_high',
]
```

Homework Request #1

```
In [51]: # Standardize the predictor variables
scaler = StandardScaler()
X = pd.DataFrame(scaler.fit_transform(X), columns=list_predictors)

In [52]: # Split data into training set
X_Train = df_loan.loc(df_loan['partition'] == 'train')[list_predictors]
y_Train = df_loan.loc(df_loan['partition'] == 'train')['loan_status']
X_Test = df_loan.loc(df_loan['partition'] == 'test')[list_predictors]
y_Test = df_loan.loc(df_loan['partition'] == 'test')['loan_status']
```

Homework Question 2

```
In [53]: from sklearn.linear_model import LogisticRegression
list_LogisticRegression = [
    LogisticRegression(
        penalty = penalty,
        solver = 'lbfgs',
        l1_ratio = 0.5, # ignored for non-elasticnet
        max_iter = 10000
    ) for penalty in (None, 'l2', 'l1', 'elasticnet')
]
list_fit = [model.fit(
    X = X_Train,
    y = y_Train
) for model in list_LogisticRegression]

d:\python\lib\site-packages\sklearn\linear_model\_logistic.py:1197: UserWarning: l1_ratio parameter is only used when penalty is 'elasticnet'. Got [penalty=None]
warnings.warn(
d:\python\lib\site-packages\sklearn\linear_model\_logistic.py:1197: UserWarning: l1_ratio parameter is only used when penalty is 'elasticnet'. Got [penalty=l2]
warnings.warn(
d:\python\lib\site-packages\sklearn\linear_model\_logistic.py:1197: UserWarning: l1_ratio parameter is only used when penalty is 'elasticnet'. Got [penalty=l1]
warnings.warn(

In [54]: [fit.intercept_.round(1) for fit in list_fit]
Out[54]: [array([-2.4]), array([-2.4]), array([-2.4]), array([-2.4])]

In [55]: [fit.coef_.round(1) for fit in list_fit]
Out[55]: [array([[[-0.1, -0. , 1.4, 0.4, 0.7]]],
array([[[-0.1, -0. , 1.4, 0.4, 0.7]]],
array([[[-0.1, -0. , 1.4, 0.4, 0.7]]],
array([[[-0.1, -0. , 1.4, 0.4, 0.7]])]

In [56]: list_predict = [fit.predict(X = X[list_predictors]) for fit in list_fit]
list_predict_proba = [fit.predict_proba(X = X[list_predictors])[:,1] for fit in list_fit]
df_predict_proba = pd.DataFrame(list_predict_proba)
df_predict_proba
```

	0	1	2	3
0	0.287672	0.287354	0.289566	0.286647
1	0.008033	0.008113	0.008142	0.008127
2	0.037255	0.037442	0.037468	0.037447
3	0.409555	0.409912	0.409034	0.409440
4	0.418099	0.417833	0.416483	0.417090
...
44995	0.418099	0.417833	0.416483	0.417090
44996	0.287672	0.287354	0.289566	0.286647
44997	0.418099	0.417833	0.416483	0.417090
44998	0.126763	0.127174	0.127643	0.127415
44999	0.373381	0.373259	0.373703	0.373400

45000 rows x 4 columns

```
In [57]: from sklearn.metrics import roc_auc_score, roc_curve
df_loan['partition']
["train", "validation", "test"]
list_roc_auc_score = [(k, roc_auc_score(
    y_true = df_loan.loc(df_loan['partition'] == k, 'loan_status'),
    y_score = df_predict_proba.loc(df_loan['partition'] == k, 1)
) for k in range(len(list_predict_proba)) for k in ("train", "validation", "test")]
```

```
Out[57]:
```

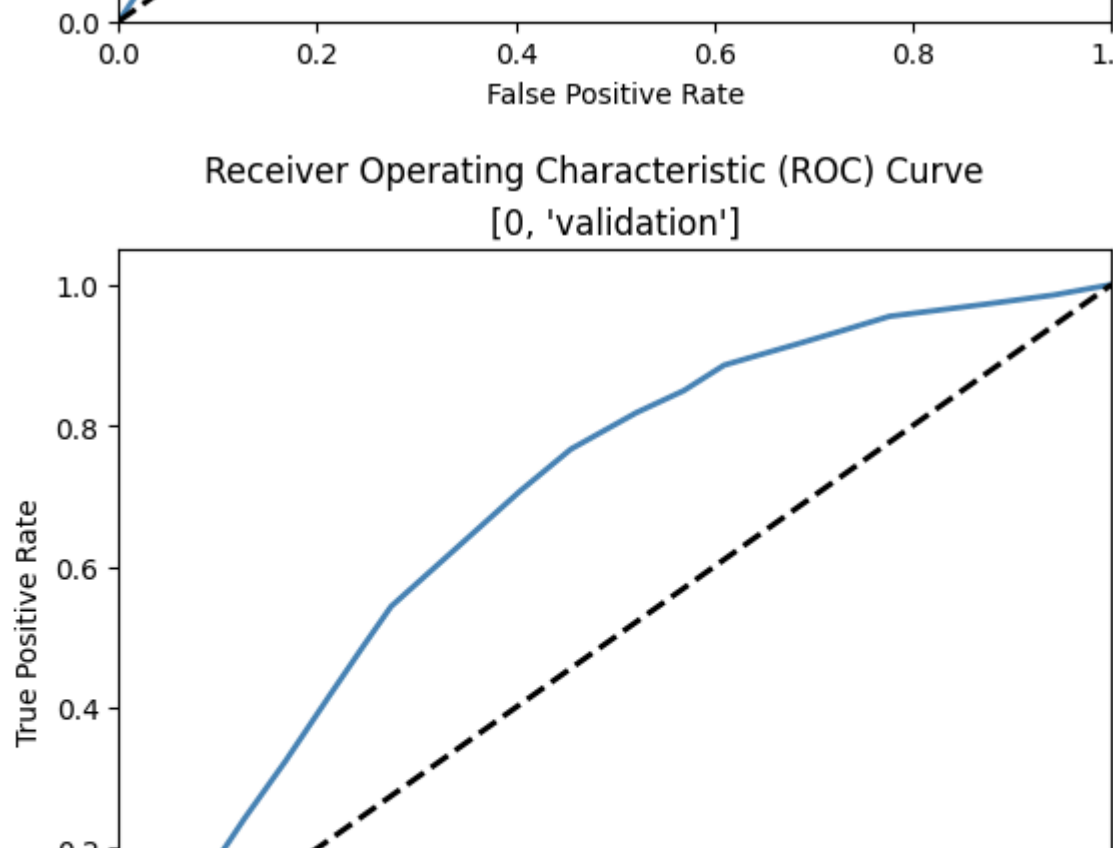
	0	1	2
0	train	0.696780	
1	validation	0.690306	
2	test	0.690289	
3	train	0.696780	
4	validation	0.690306	
5	test	0.690289	
6	train	0.696780	
7	validation	0.690306	
8	test	0.690289	
9	train	0.696780	
10	validation	0.690306	
11	test	0.690289	

```
In [58]: list_roc_curve = [(k, roc_curve(
    y_true = df_loan.loc(df_loan['partition'] == k, 'loan_status'),
    y_score = df_predict_proba.loc(df_loan['partition'] == k, 1)
) for k in range(len(list_predict_proba)) for k in ("train", "validation", "test")]
```

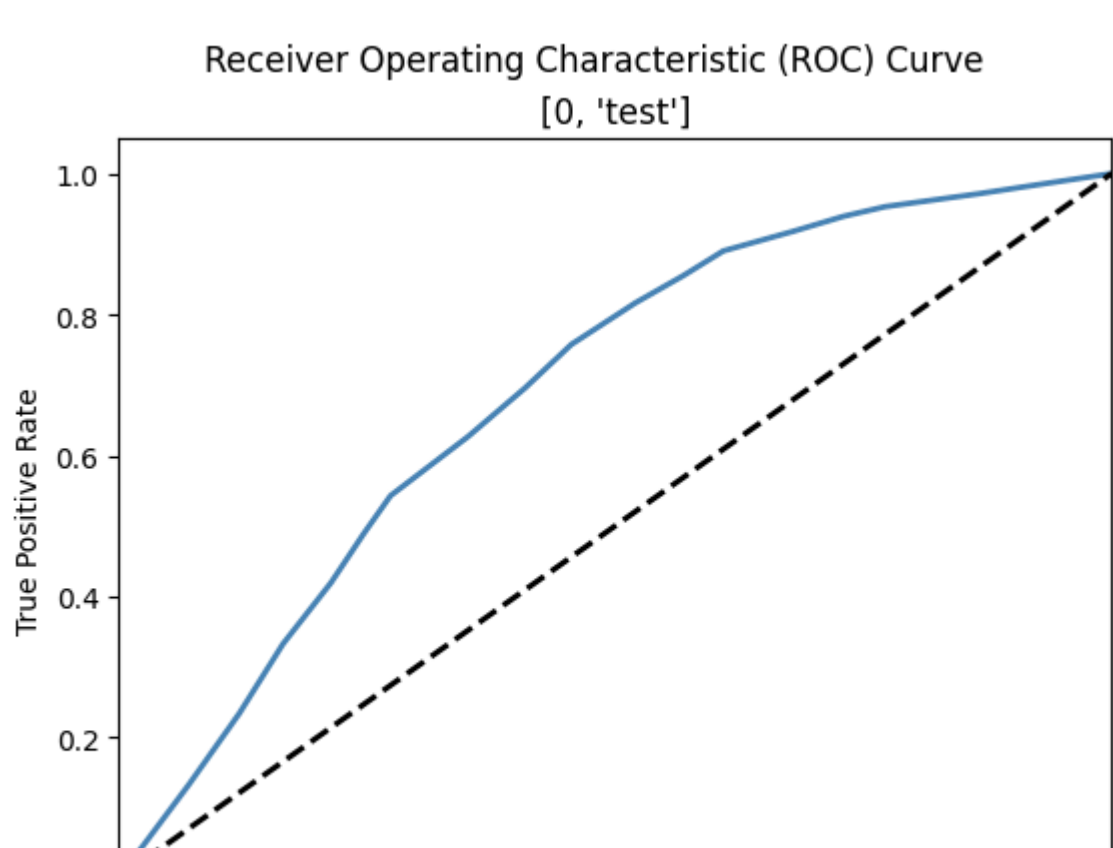
```
In [59]: import matplotlib.pyplot as plt
for j in range(len(list_roc_curve)):
    plt.figure()
    color = ['steelblue']
    for i, color in zip(range(2), colors):
        plt.plot(
            list_roc_curve[j][2][0],
            list_roc_curve[j][2][1],
            color = color,
            lw = 2
        )
        plt.plot([0, 1], [0, 1], 'k--', lw=2)
        plt.xlim(0.0, 1.0)
        plt.ylim(0.0, 1.0)
        plt.xlabel('False Positive Rate')
        plt.ylabel('True Positive Rate')
        plt.title('Receiver Operating Characteristic (ROC) Curve %i' % j)
        plt.show()
```

#fpr, tpr, thresholds = metrics.roc_curve(y, scores, pos_label=2)

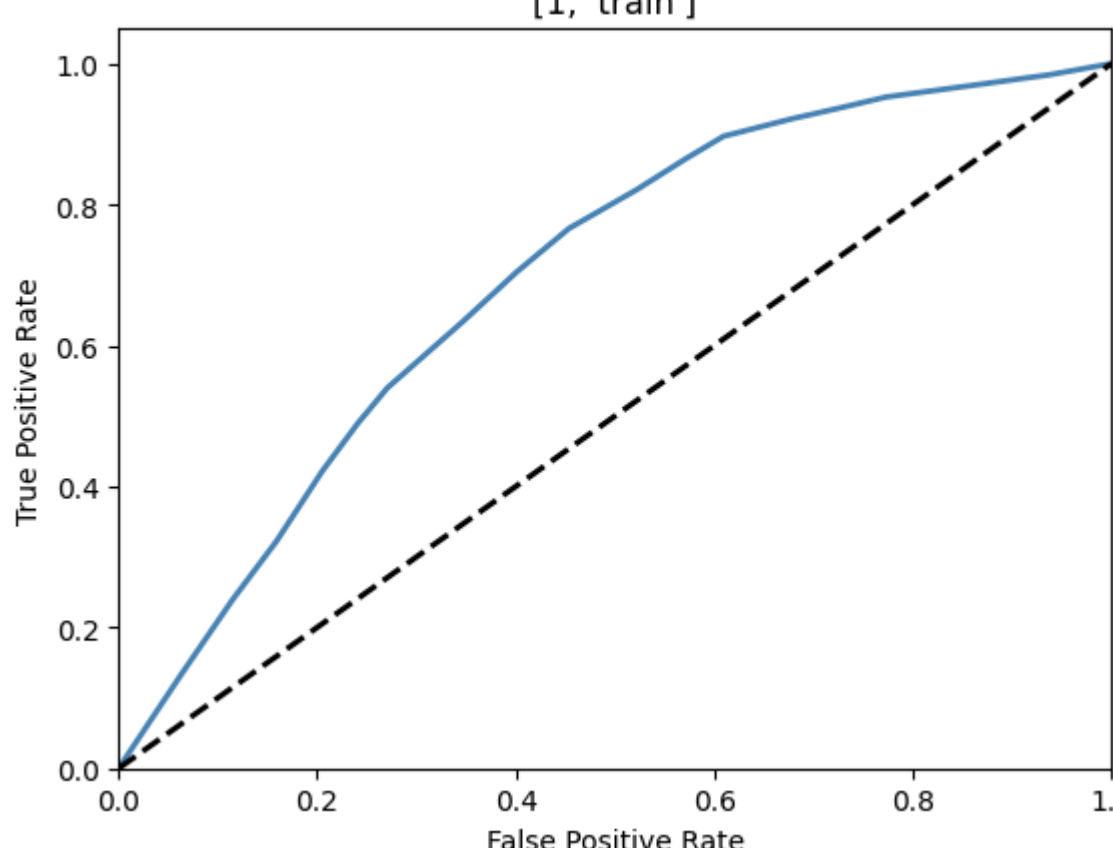
Receiver Operating Characteristic (ROC) Curve [0, 'train']



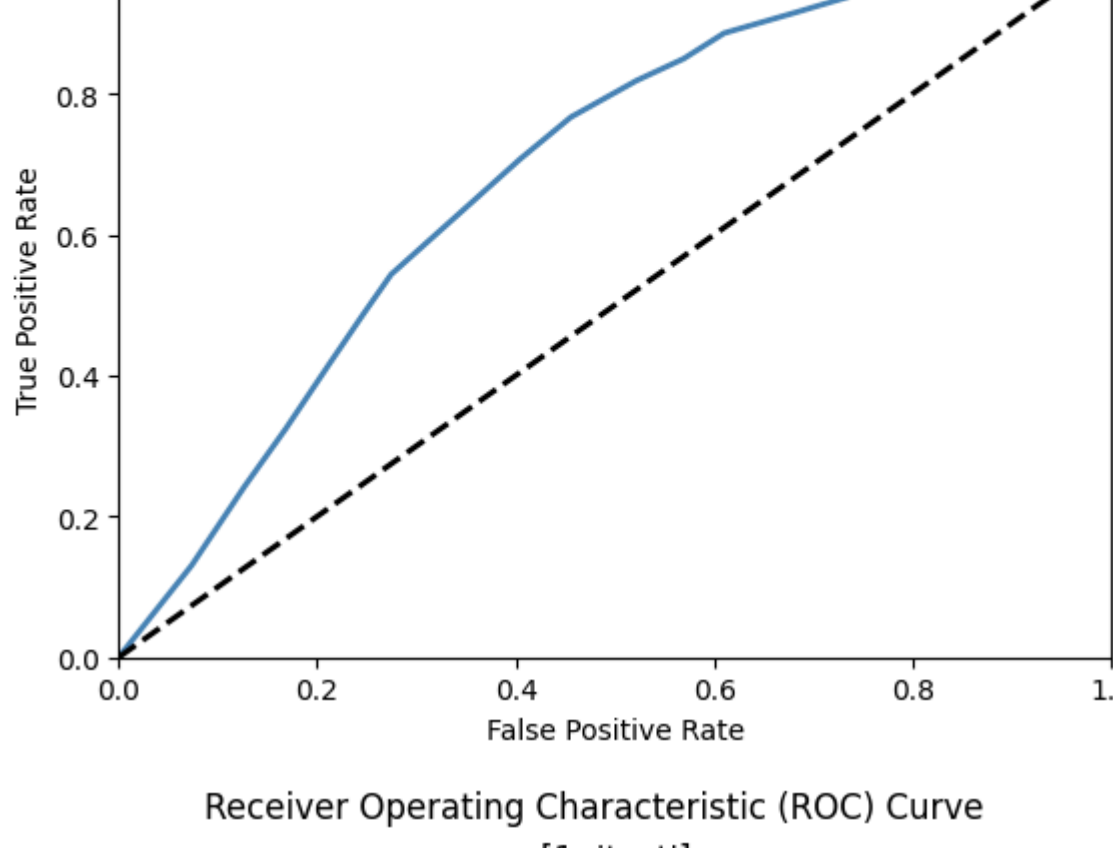
Receiver Operating Characteristic (ROC) Curve [0, 'validation']



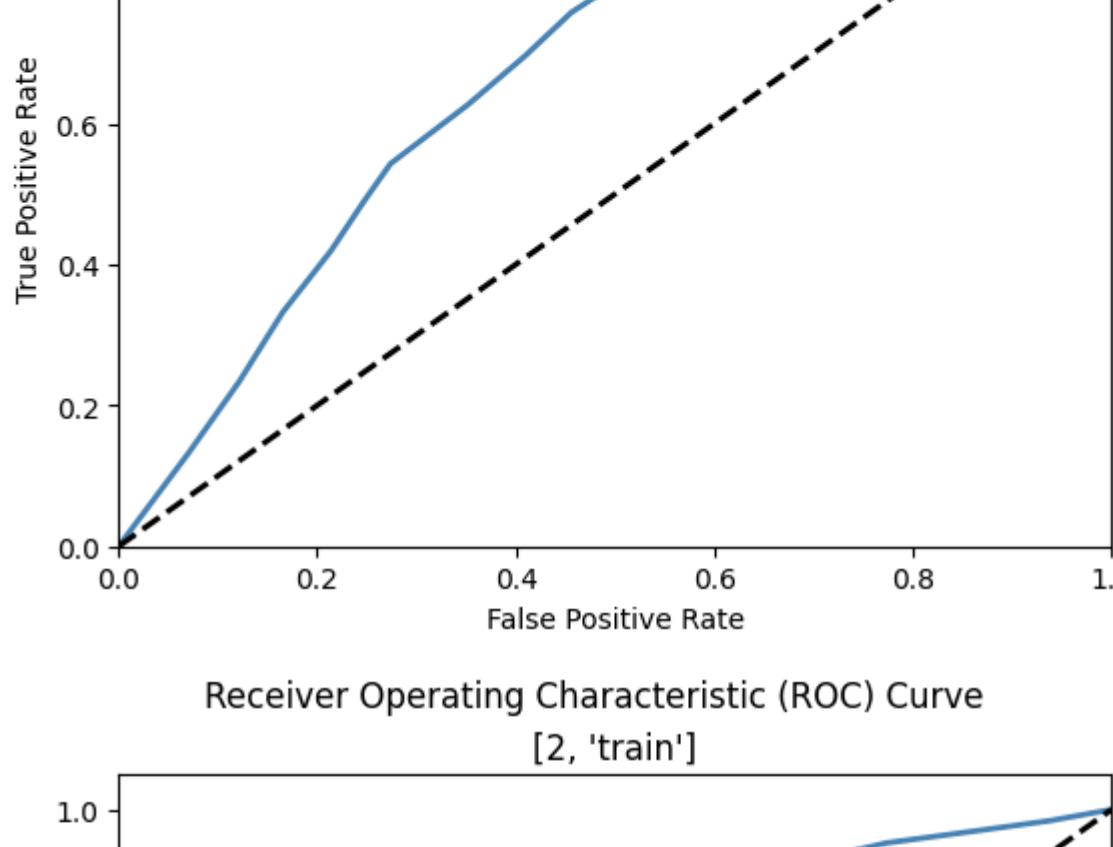
Receiver Operating Characteristic (ROC) Curve [0, 'test']



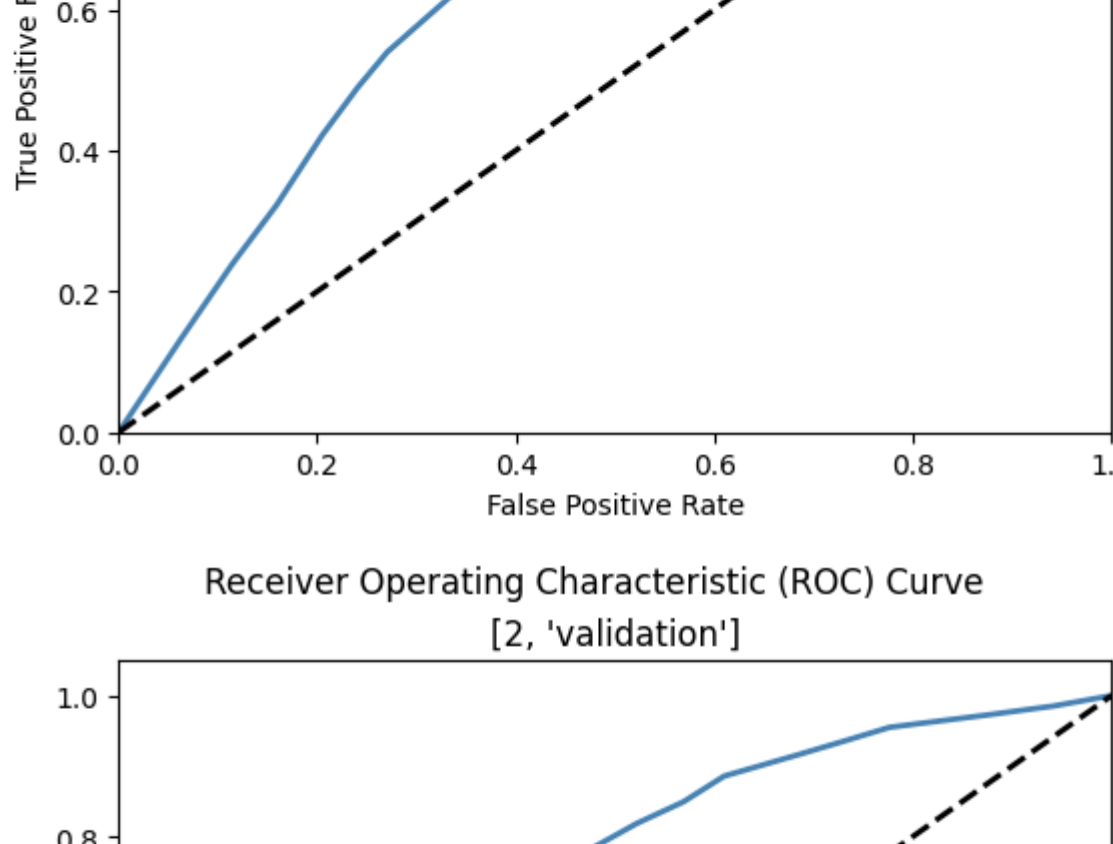
Receiver Operating Characteristic (ROC) Curve [1, 'train']



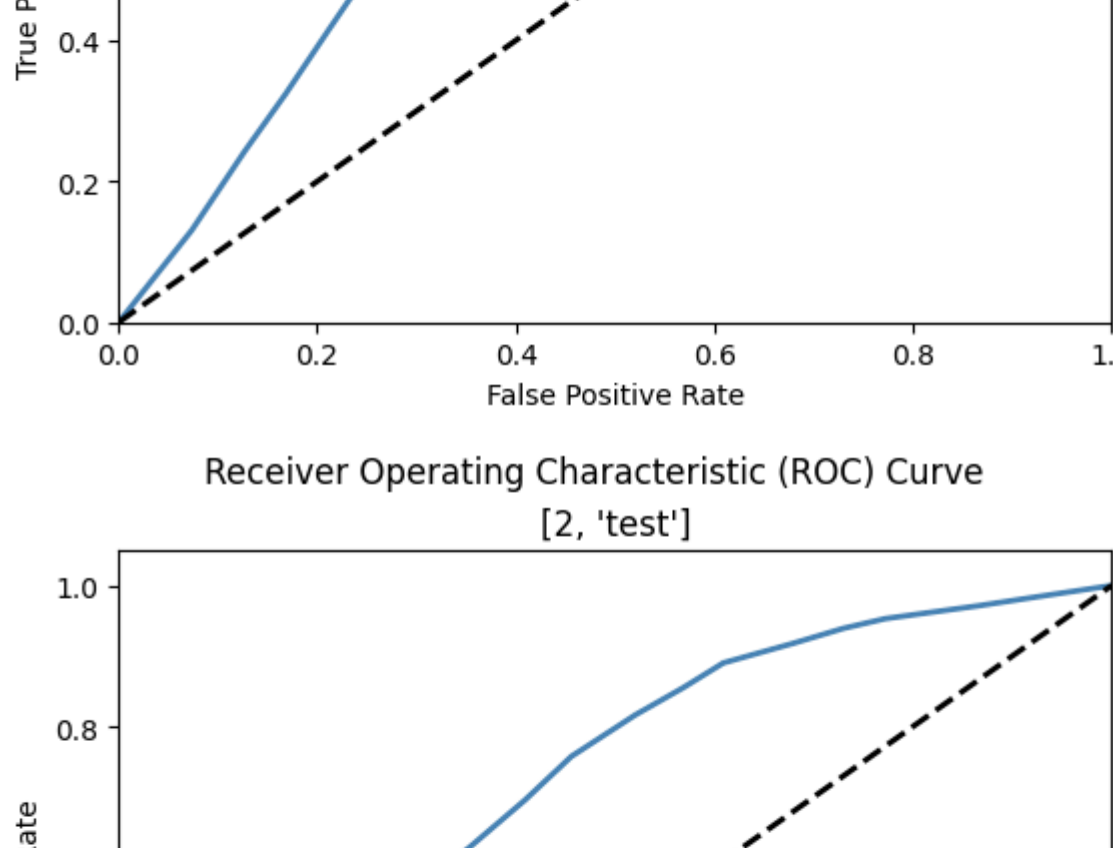
Receiver Operating Characteristic (ROC) Curve [1, 'validation']



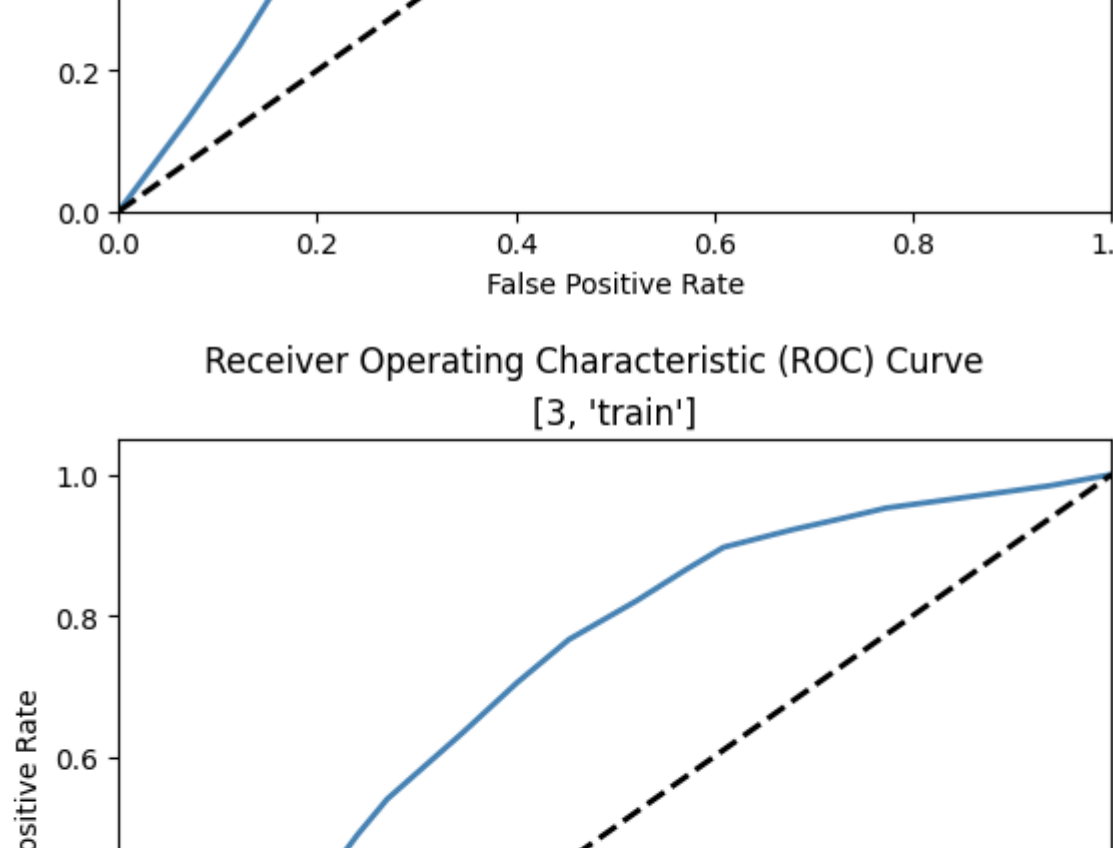
Receiver Operating Characteristic (ROC) Curve [1, 'test']



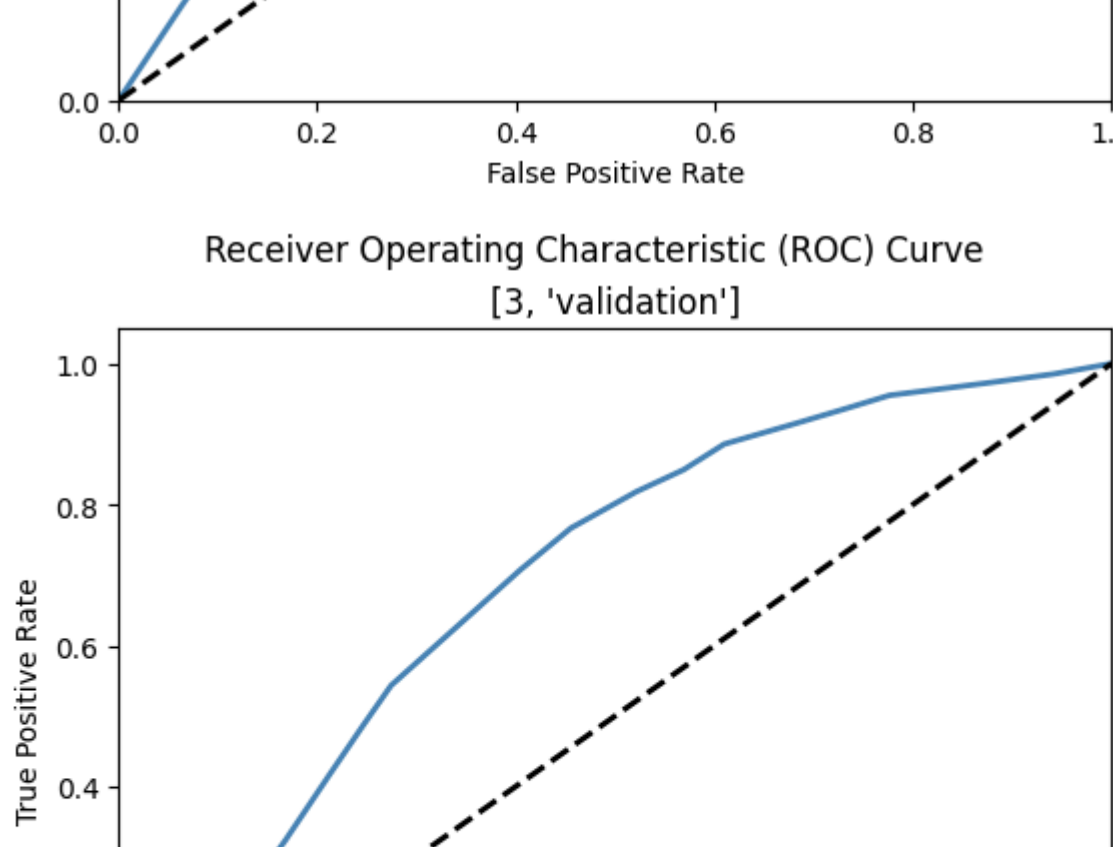
Receiver Operating Characteristic (ROC) Curve [2, 'train']



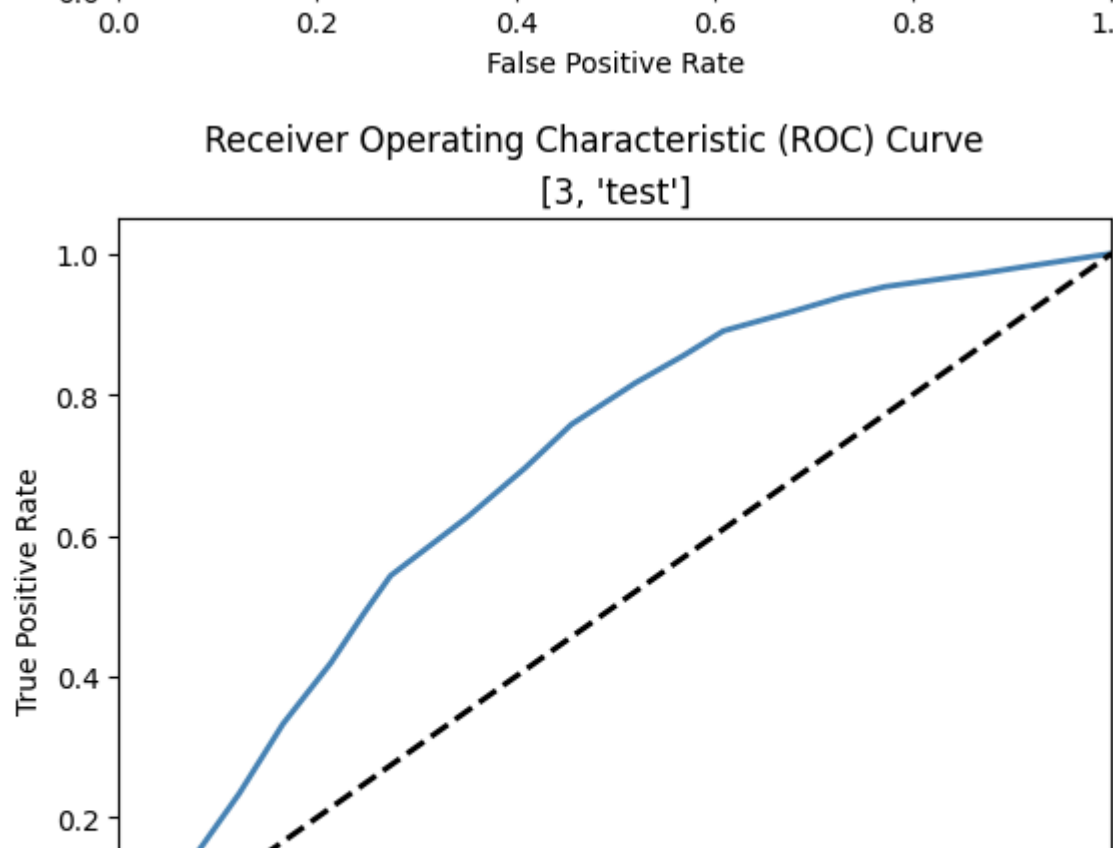
Receiver Operating Characteristic (ROC) Curve [2, 'validation']



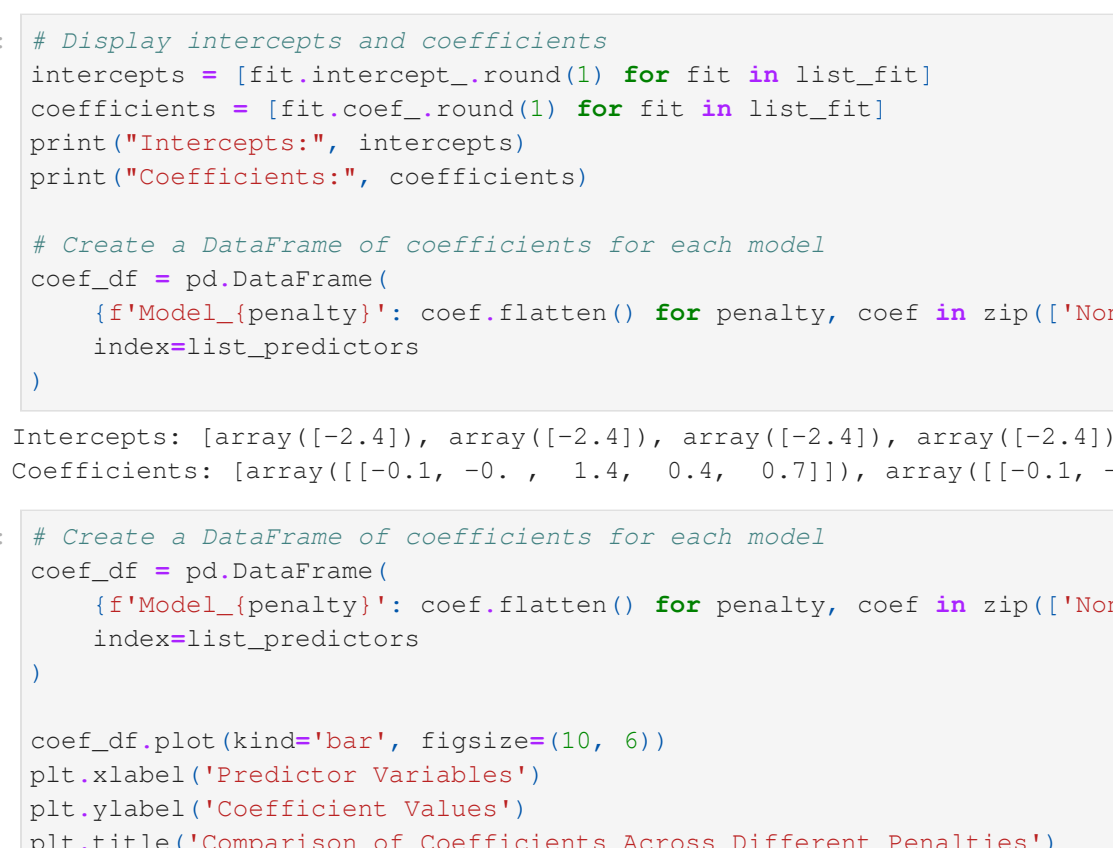
Receiver Operating Characteristic (ROC) Curve [2, 'test']



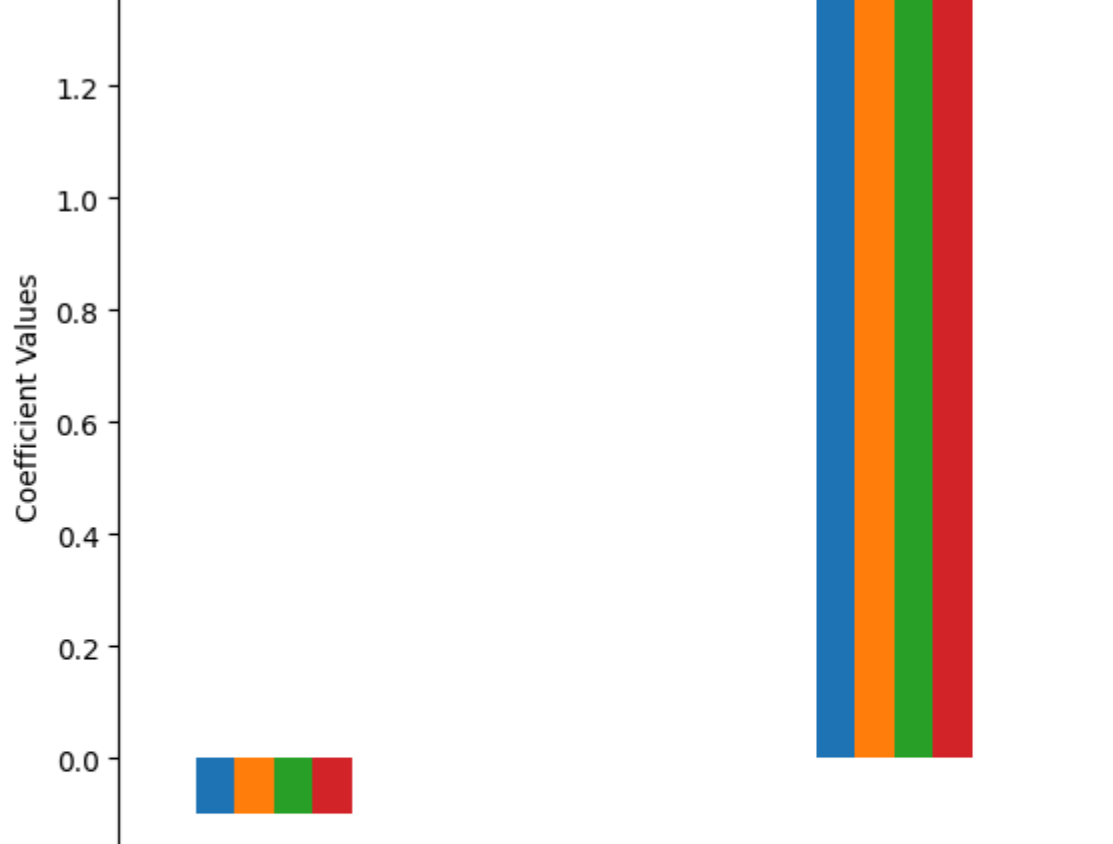
Receiver Operating Characteristic (ROC) Curve [3, 'train']



Receiver Operating Characteristic (ROC) Curve [3, 'validation']



Receiver Operating Characteristic (ROC) Curve [3, 'test']



```
In [60]: # Display intercepts and coefficients
intercepts = [fit.intercept_.round(1) for fit in list_fit]
coefficients = [fit.coef_.round(1) for fit in list_fit]
print("Intercepts:", intercepts)
print("Coefficients:", coefficients)

# Create a DataFrame of coefficients for each model
coef_df = pd.DataFrame(
    [{"model_penalty": coef.flatten() for penalty, coef in zip(['None', 'l2', 'l1', 'elasticnet'], coefficients)},
    index=list_predictors
]

Intercept: [array([-2.4]), array([-2.4]), array([-2.4]), array([-2.4])]
Coefficients: [array([[[-0.1, -0. , 1.4, 0.4, 0.7]]], array([[[-0.1, -0. , 1.4, 0.4, 0.7]]], array([[[-0.1, -0. , 1.4, 0.4, 0.7]]], array([[[-0.1, -0. , 1.4, 0.4, 0.7]])]

In [61]: # Create a DataFrame of coefficients for each model
coef_df = pd.DataFrame(
    [{"model_penalty": coef.flatten() for penalty, coef in zip(['None', 'l2', 'l1', 'elasticnet'], coefficients)},
    index=list_predictors
]

coef_df.plot(kind='bar', figsize=(10, 6))
plt.xlabel('Predictor Variables')
plt.ylabel('Coefficient Values')
plt.title('Comparison of Coefficients Across Different Penalties')
plt.show()
```



Homework Quesiton 3 NOTES ABOUT DIFFERENCE

It looks like model l2 and l1 have penalties that disallow overfitting. elasticnet has bot l1 and l2 in it. lastly, l1 is able to make the model ignore features, making the coefficient 0.