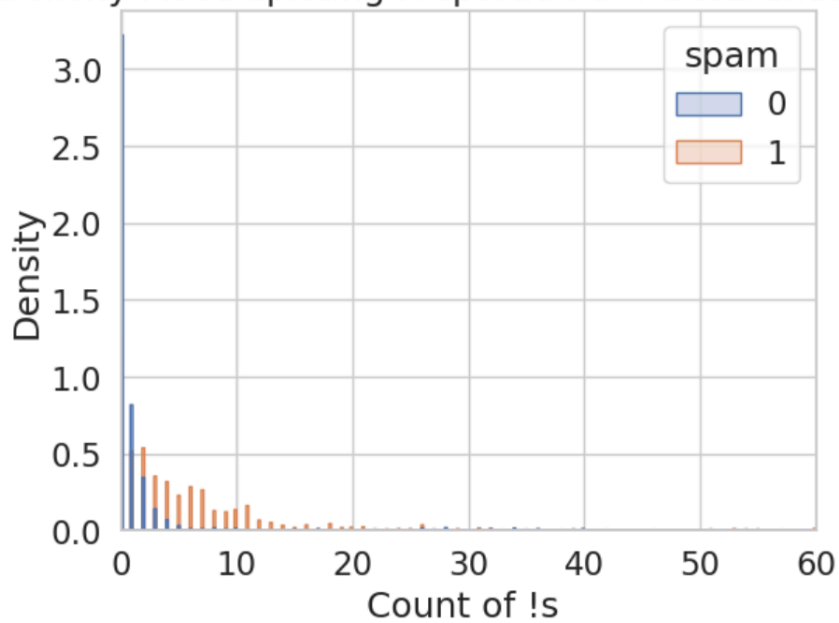
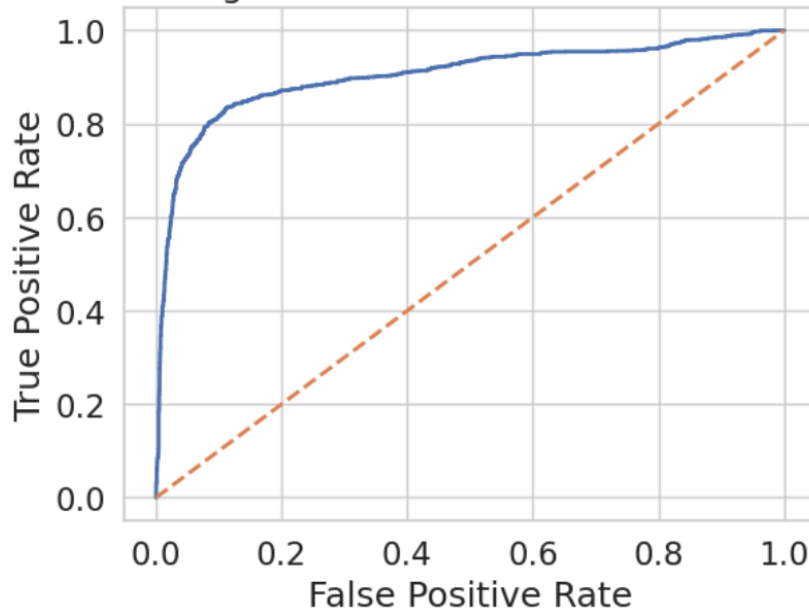


Histogram Density Plot Depicting Proportion of ! Occurances in Spam, Ham



ROC Curve: Assessing Model Performance at Various Threshold Levels



```
# Define your processing function, processed data, and model here.
# You may find it helpful to look through the rest of the questions first!
def processingFunction(data, spamWords):
    for word in spamWords:
        data[word + '_count'] = 0
        data[word + '_count'] = data['email'].apply(lambda x: countWords(x, word))

    data["!count"] = data["email"].str.count("!")
    data["emailWordCt"] = data["email"].apply(len)

    for word in ["spam", "id", "subject", "email"]:
        if word in data.columns:
            data = data.drop(columns = [word])

    return data
def countWords(text, word):
    return text.count(word)
someWords= ["free","win","selected","winner", "credit","click","urgent", "100%",
            "#1", "trial", "earn", "guarantee", "viagra", "dick","ad ", "scam",
            "income","url", "div", "font", "body", "link",
            "html", "http", "head", "hot", "debt",
            "prize", "$","\n", "=", "type"]
myModel = LogisticRegression(C=.3, max_iter=10000, solver="liblinear")

X_train = processingFunction(train,someWords)
Y_train = np.array(train["spam"])
myModel.fit(X_train,Y_train)
```

```
LogisticRegression
LogisticRegression(C=0.3, max_iter=10000, solver='liblinear')
```