# Module 3 Summary

October 13, 2019

Paul Raymond
February 2019 cohort
Brandon Lewis

## 0.1 Project focus

My primary focus for my Module 3 Project will be executing various classification methods on a set of data from the New York Times. Additionally, I will integrate Time Series and NLP analysis into the workflow.

I've chosen a dataset provided by The New York Times that includes data on articles and comments written during Jan-April 2017 and Jan-April 2018.

I plan to measure whether features - numerical, categorical, text - can be used to accurately predict which comments are selected by NYTimes editors.

## 0.2 Exploratory Data Analysis

I'll be working with the following files made available by the New York Times on Kaggle: **Articles and Comments submitted duirng January - April 2017 and 2018** After combining and processing the datasets, here's a breakdown on the Articles and Comments data:

*Breakdown by EditorsSelection*

*Numerical features of focus*

There isn't consistent dilinearation between selected and nonselected comments with respect to numerical features. The word count distribution for nonselected comments is skewed right whereas that for selected comments is more normally distribution.

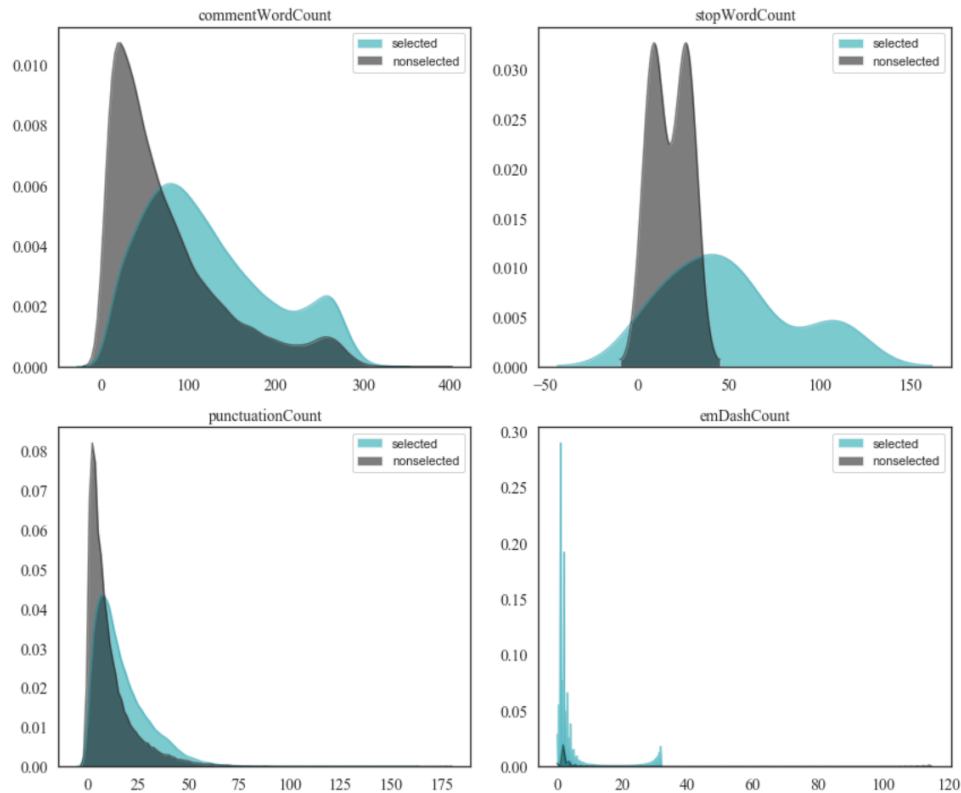The violin plots confirm this as well.

*Data breakdown*

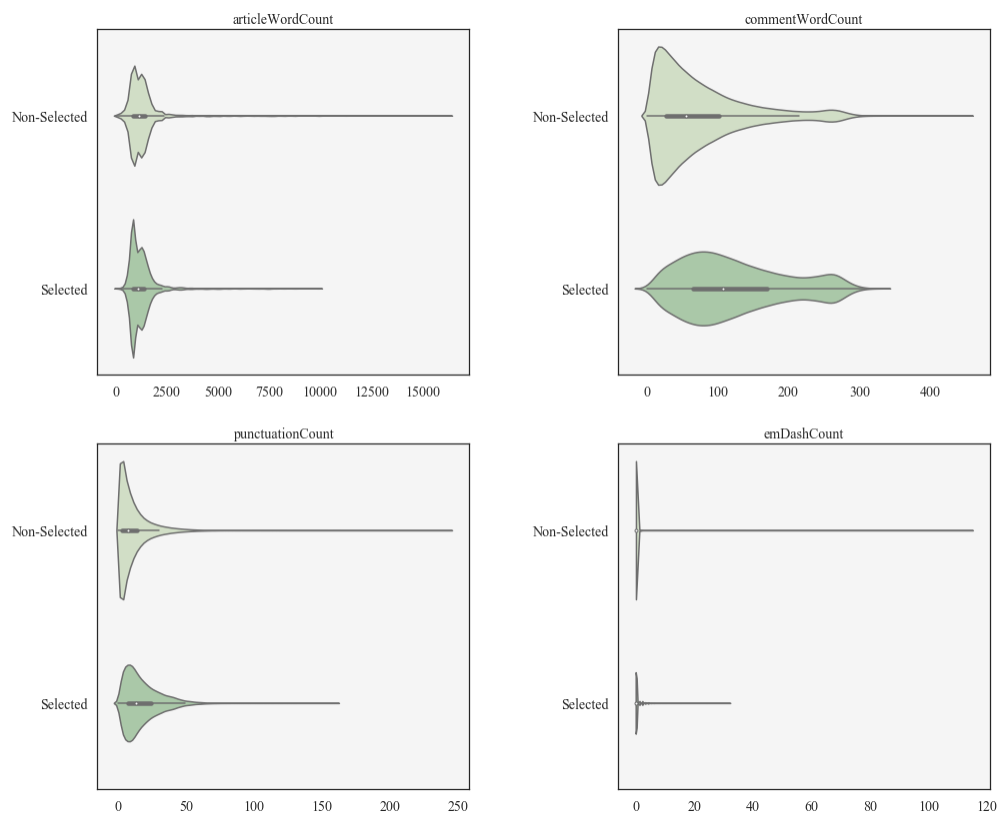*Type of material (Total articles & comments)*

## Some data and meta-comments

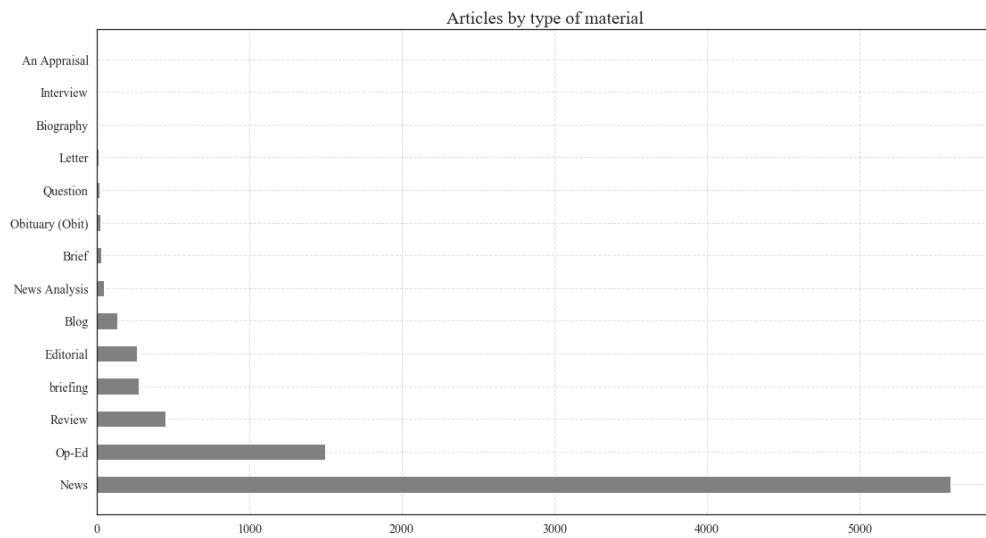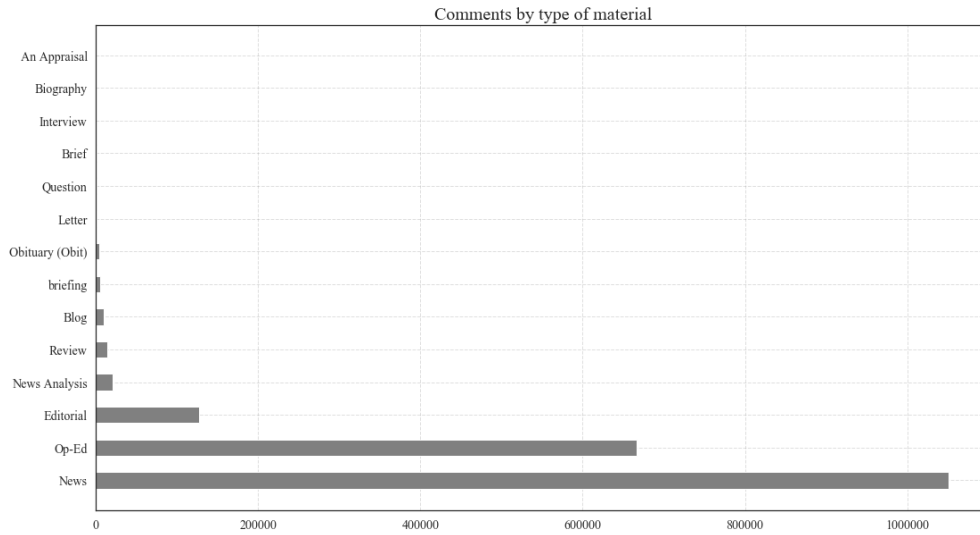| | Articles | Comments |
|---|---|---|
| | **during Jan – April (2017 & 2018)** | |
| **2 Datasets:** | • **8,339 articles** written<br>• **17 features** for each article | • **1,899,975 comments** written<br>• **38 features** for each comment |
| **Range:** | • **Longest article:** "On the shooting in Florida, Student Activism.." **(16,336 words)**<br>• **Shortest article:** Unknown **(11 words)**<br>• **Most articles: Deb Amlen**<br>• **Most words: Editorial Board** | • **Longest comment:** "Beautifully written and deeply touching. It makes me wonder how many of us also have college classmates who died while homeless and suffering from mental illness.." **(452 words)**<br>• **Shortest comment: (4232 with one word)** |
| | • **News** was type of material with **most articles (5,596)**<br>• **OpEd** was category with **most articles (1,528)** | • **News** was type of material with **most comments (1M+)**<br>• **OpEd** was category with **most comments (671893)** |

title



title

"Plot"



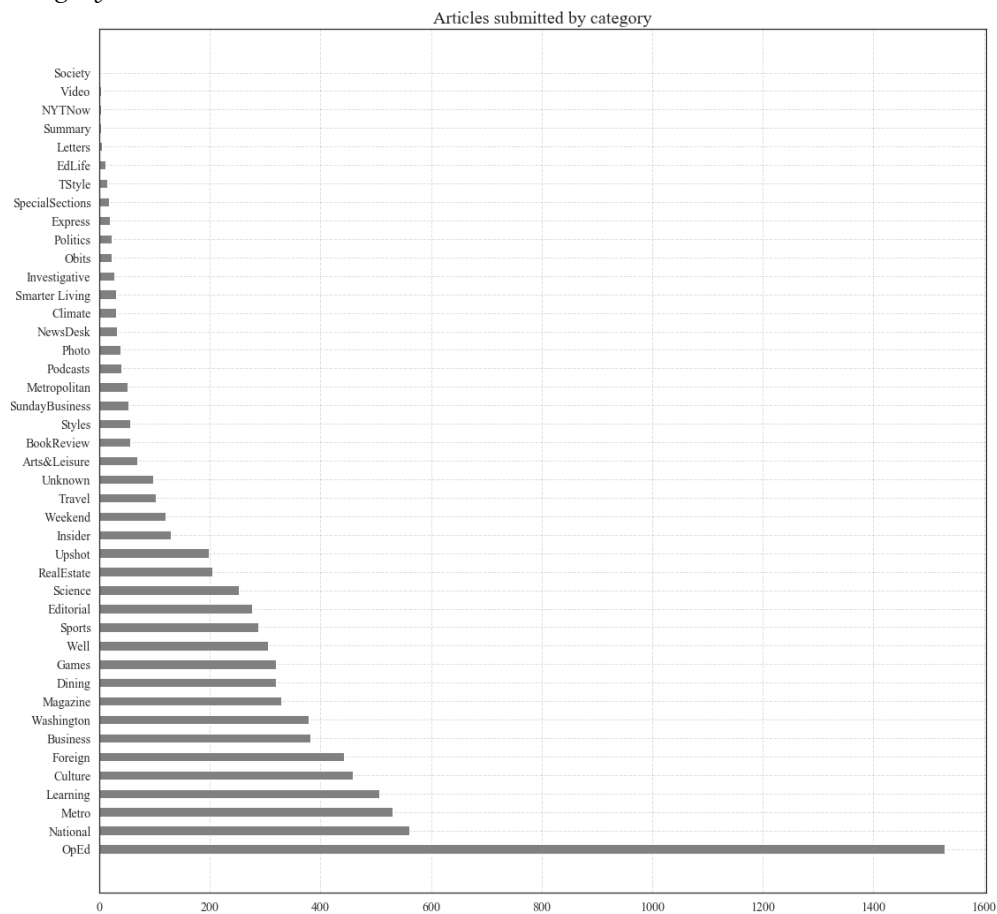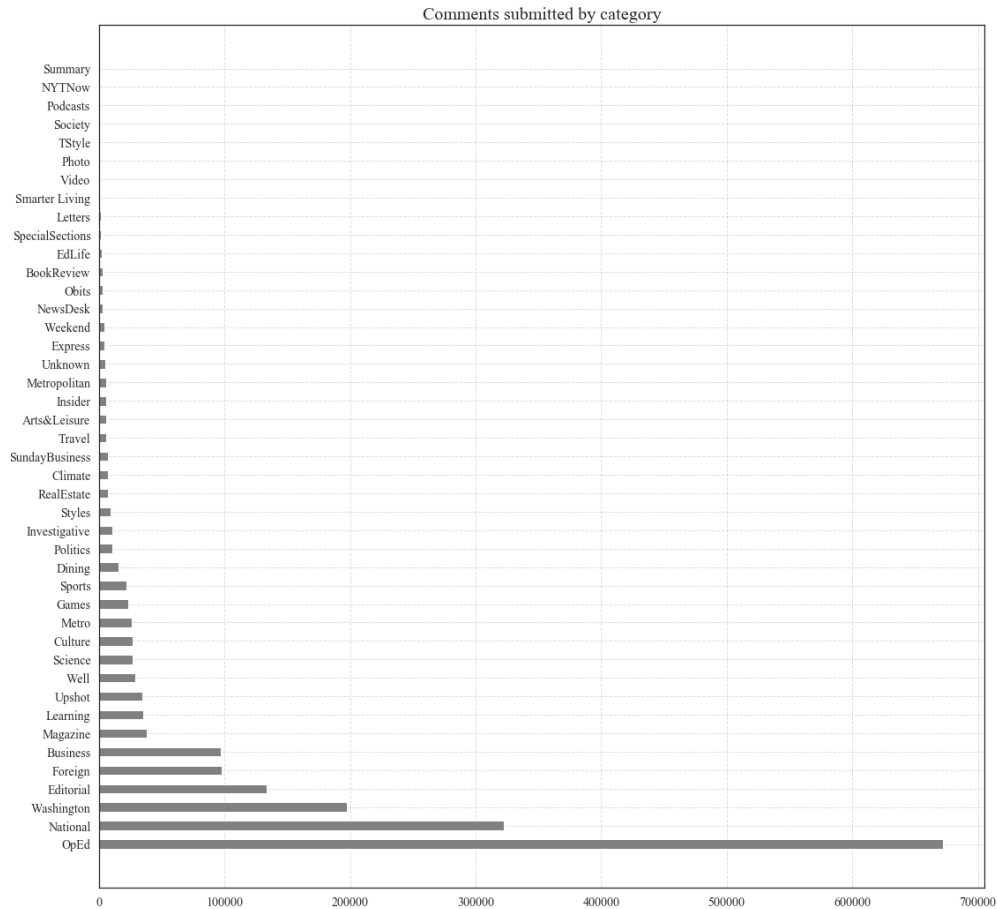title

Comments by type of material
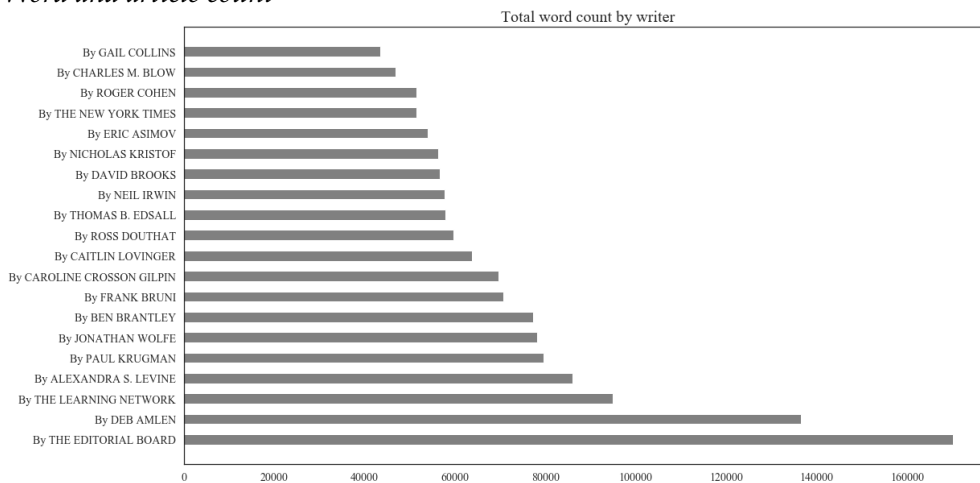
There appears to be more comments submitted per Op-Eds than ratio for any other main category — what I would expect.
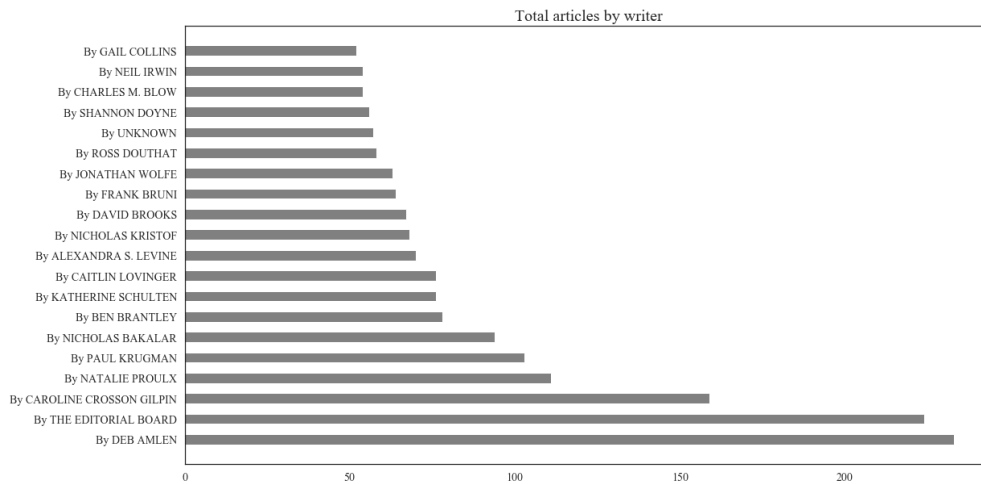
*Category (Total articles & comments)*



Articles submitted by category

Comments submitted by category

There appears to be more comments submitted per Op-Eds than ratio for any other main category — what I would expect.

*Word and article count*
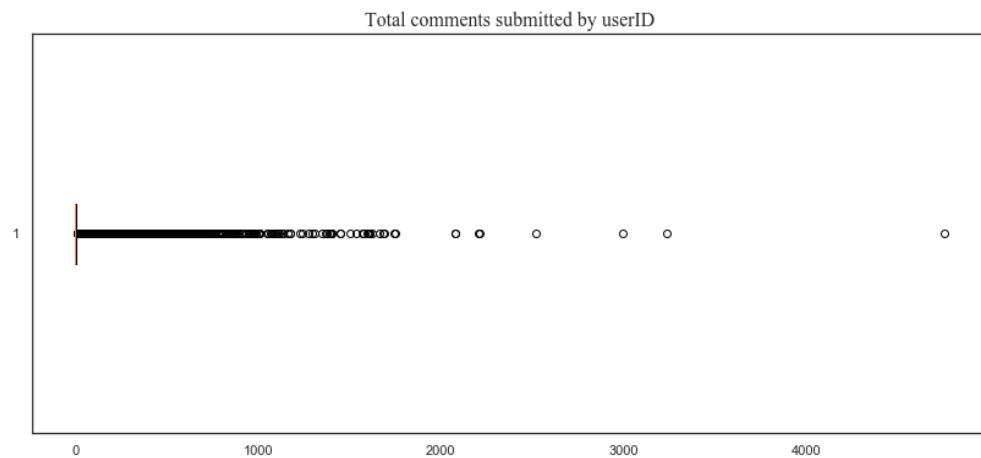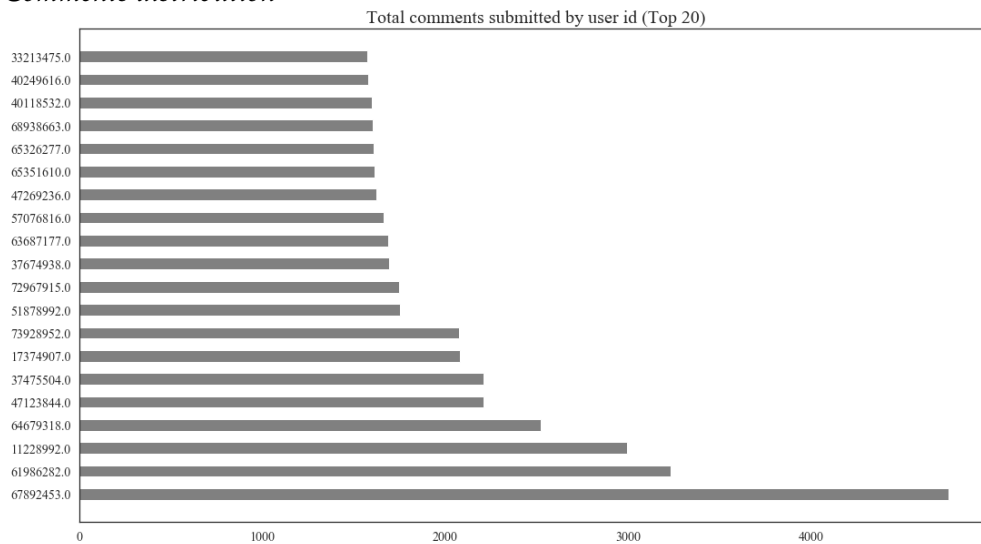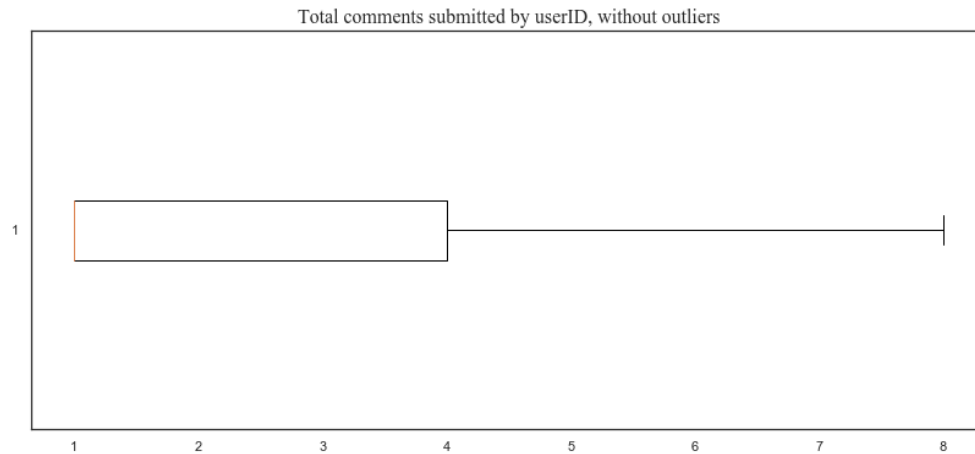


Total word count by writer

**Total articles by writer**

| Writer | |
|---|---|
| By GAIL COLLINS | |
| By NEIL IRWIN | |
| By CHARLES M. BLOW | |
| By SHANNON DOYNE | |
| By UNKNOWN | |
| By ROSS DOUTHAT | |
| By JONATHAN WOLFE | |
| By FRANK BRUNI | |
| By DAVID BROOKS | |
| By NICHOLAS KRISTOF | |
| By ALEXANDRA S. LEVINE | |
| By CAITLIN LOVINGER | |
| By KATHERINE SCHULTEN | |
| By BEN BRANTLEY | |
| By NICHOLAS BAKALAR | |
| By PAUL KRUGMAN | |
| By NATALIE PROULX | |
| By CAROLINE CROSSON GILPIN | |
| By THE EDITORIAL BOARD | |
| By DEB AMLEN | |

```
In [ ]:
```

## *Comments distribution*

**Total comments submitted by user id (Top 20)**

| user id |
|---|
| 33213475.0 |
| 40249616.0 |
| 40118532.0 |
| 68938663.0 |
| 65326277.0 |
| 65351610.0 |
| 47269236.0 |
| 57076816.0 |
| 63687177.0 |
| 37674938.0 |
| 72967915.0 |
| 51878992.0 |
| 73928952.0 |
| 17374907.0 |
| 37475504.0 |
| 47123844.0 |
| 64679318.0 |
| 11228992.0 |
| 61986282.0 |
| 67892453.0 |

**Total comments submitted by userID**

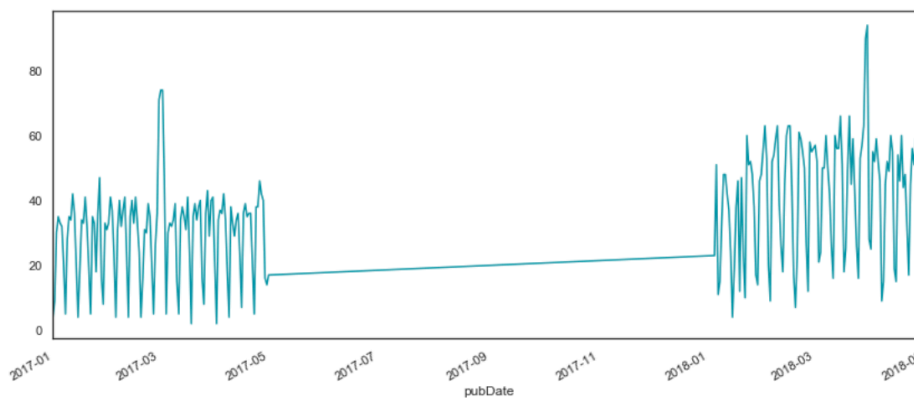Total comments submitted by userID, without outliers

As evidenced above, there are significant outliers, yet the average reader has submitted closer to 1 comment.
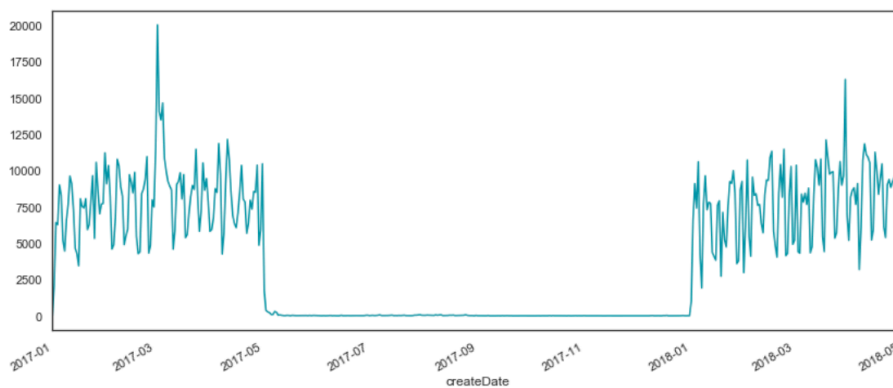
Next, I'll explore the changes in our data over the two periods: Jan - May 2017 and Jan - May 2018.

## 0.3  Time-series analysis
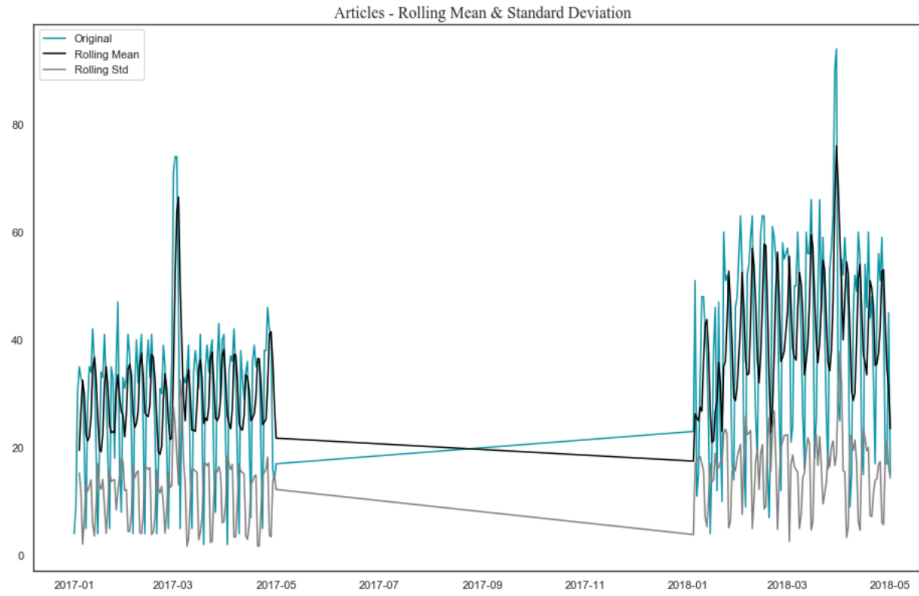
*Submissions in Jan - April (2017 and 2018)*



Visually, there appears to be an increase in articles from the period in 2017 to that in 2018.
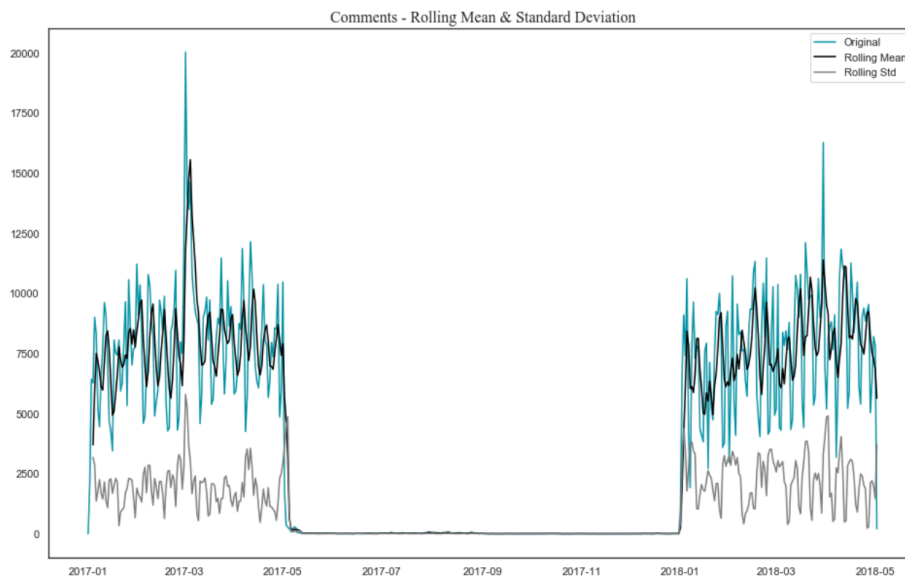


However, an increase in comments during the similar period as Articles isn't apparent from the plot of submitted comments over the two year period.

*Adjusted trends in 2017 and 2018*
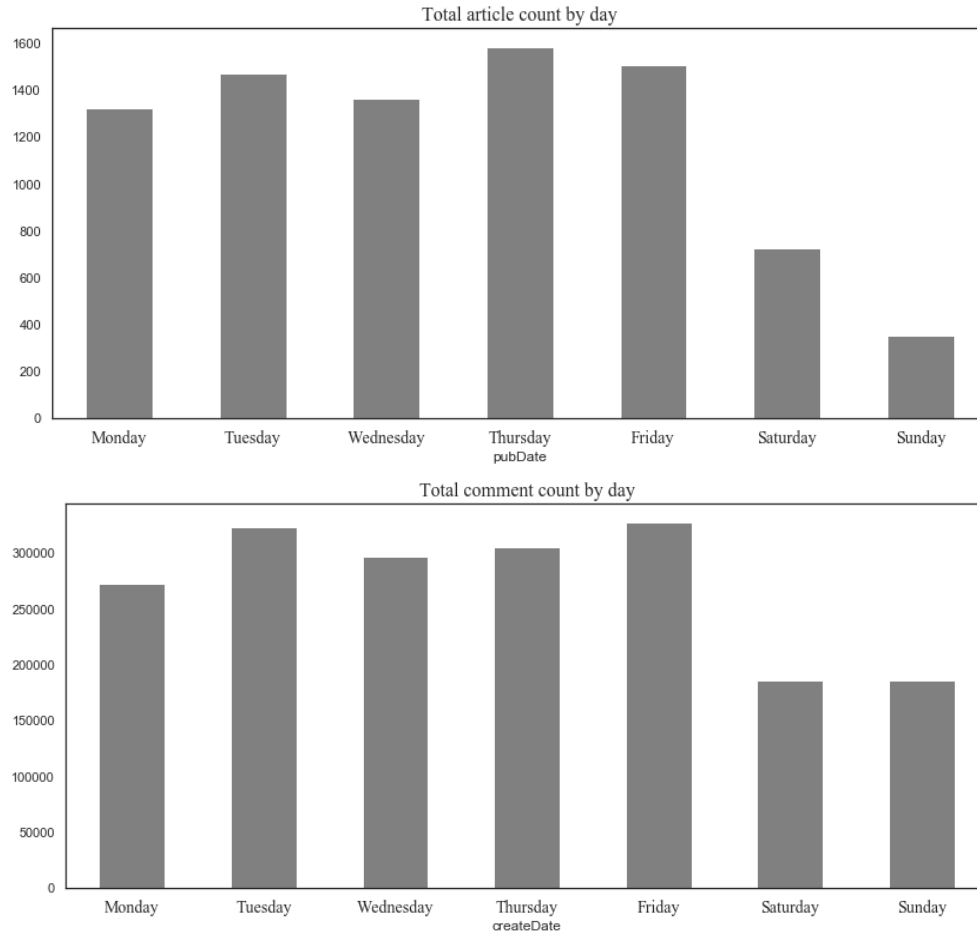
Articles - Rolling Mean & Standard Deviation

From the plot, it's apparent that the mean article count has increased and that there isn't considerable increase in variation of article counts between the two periods.



Comments - Rolling Mean & Standard Deviation

Upon further inspection, comments display consistent trends between both periods. These results are confirmed with a Dickey-Fuller test.

*Daily totals*

Total article count by day


Total comment count by day

The total articles and comments submitted by day display similar distributions. As expected, there is a drop-off in articles over the weekend, with a similar decrease in submitted comments. This is consistent with the weekly seasonality seen in the full time-series plot above.

## 0.4  Classification

I create classification models with various approaches and through two different lenses: 1.predictive performance of non-text features and 2. predictive performance of text-features. In addition to finding the best model, I'm curious to compare the performance of both approaches.
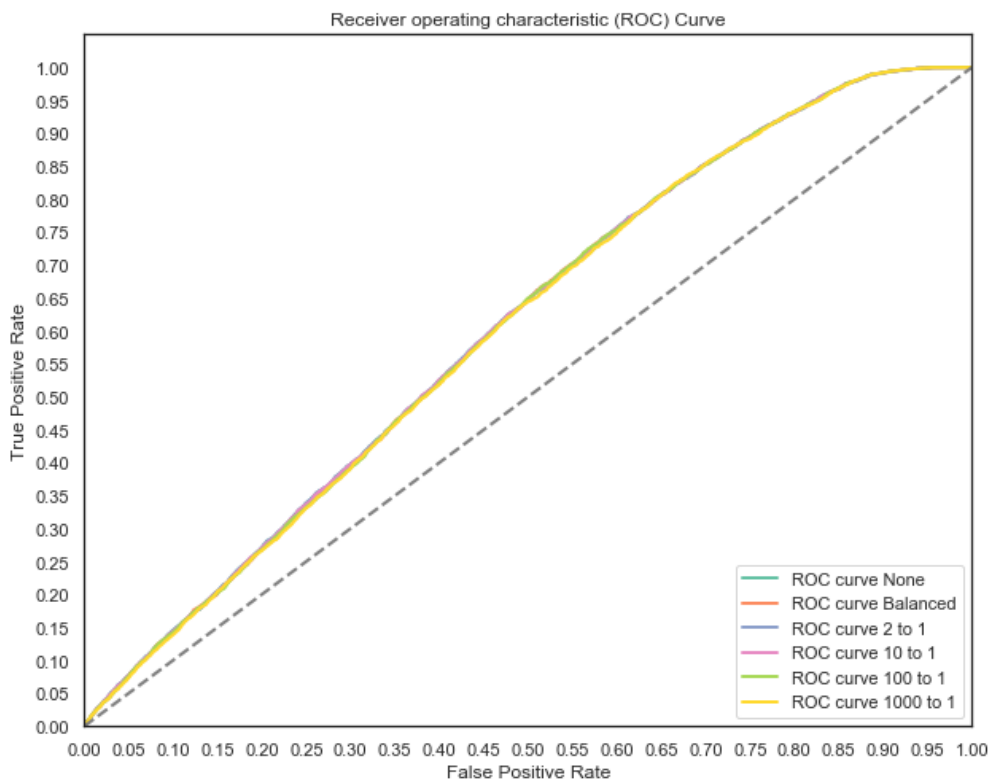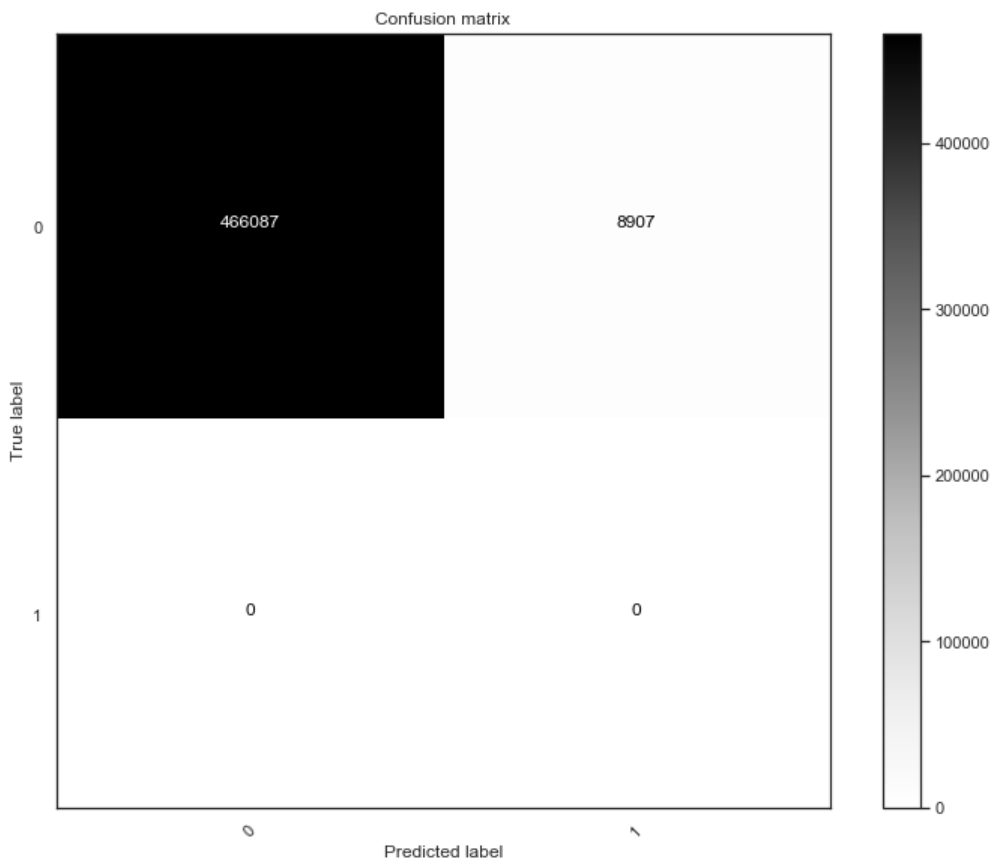
The classification model workflow is as follows: 1. **Create two datasets** - one with original columns and one that includes original and calculated columns (updated dataset) 2. **Classification models for non-text features** A. Complete baseline model (Logistic Regression) B. Complete SMOTE model (Logistic Regression) C. Complete Under-sampled model (Logistic Regression) D. Complete pipeline models (Logistic, SVM, RandomForest) E. Complete model with gridsearch 3. **Classification models for text features** A. Complete baseline model (Logistic Regression) B. Complete SMOTE model (Logistic Regression) C. Complete Under-sampled model (Logistic Regression) D. Complete pipeline models (Logistic, SVM, RandomForest) E. Complete model with gridsearch

*Classification models with solely text feature*
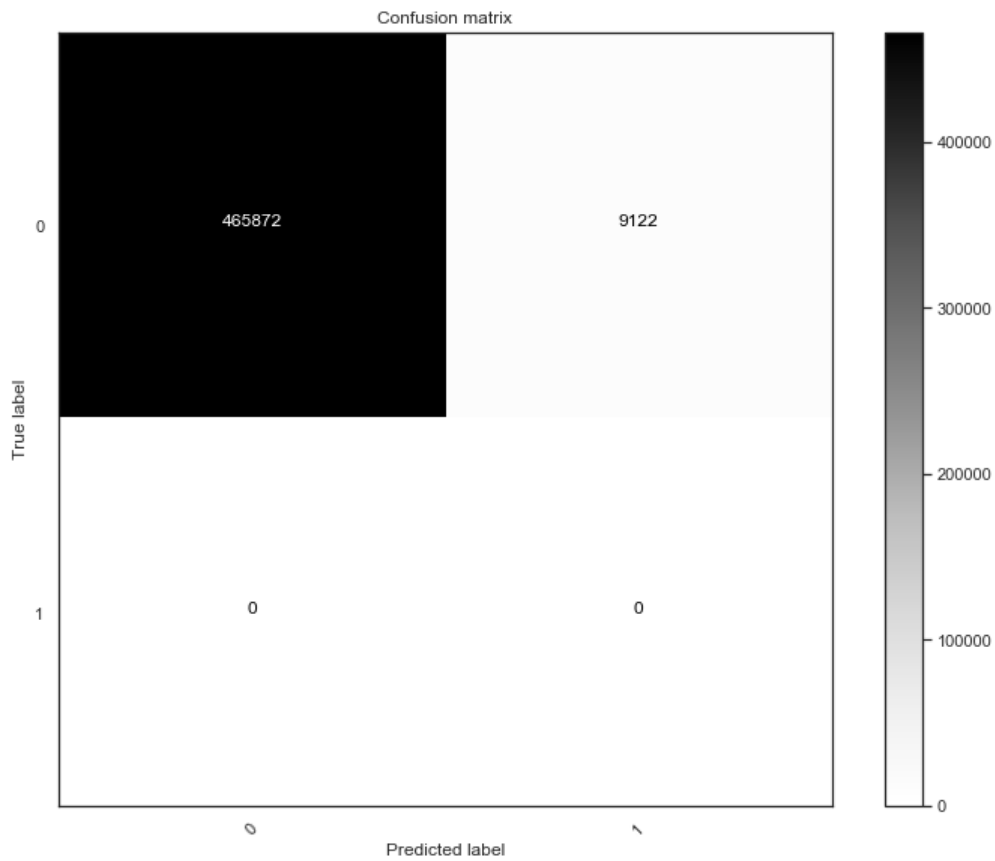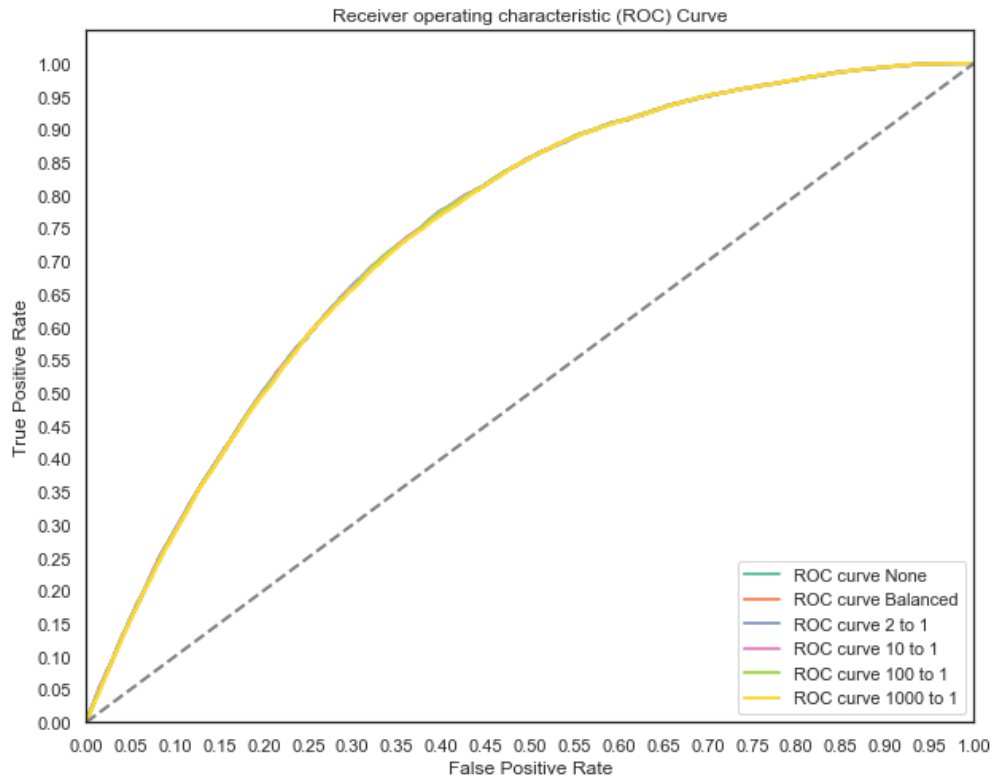**Baseline models**
*Logistic regression*

9

*Original dataset*

## Confusion matrix



## Receiver operating characteristic (ROC) Curve

As evident in the plots above, the model is affected by the class imbalance present in the target feature. As a result, ......

*Updated dataset*

Confusion matrix

Receiver operating characteristic (ROC) Curve
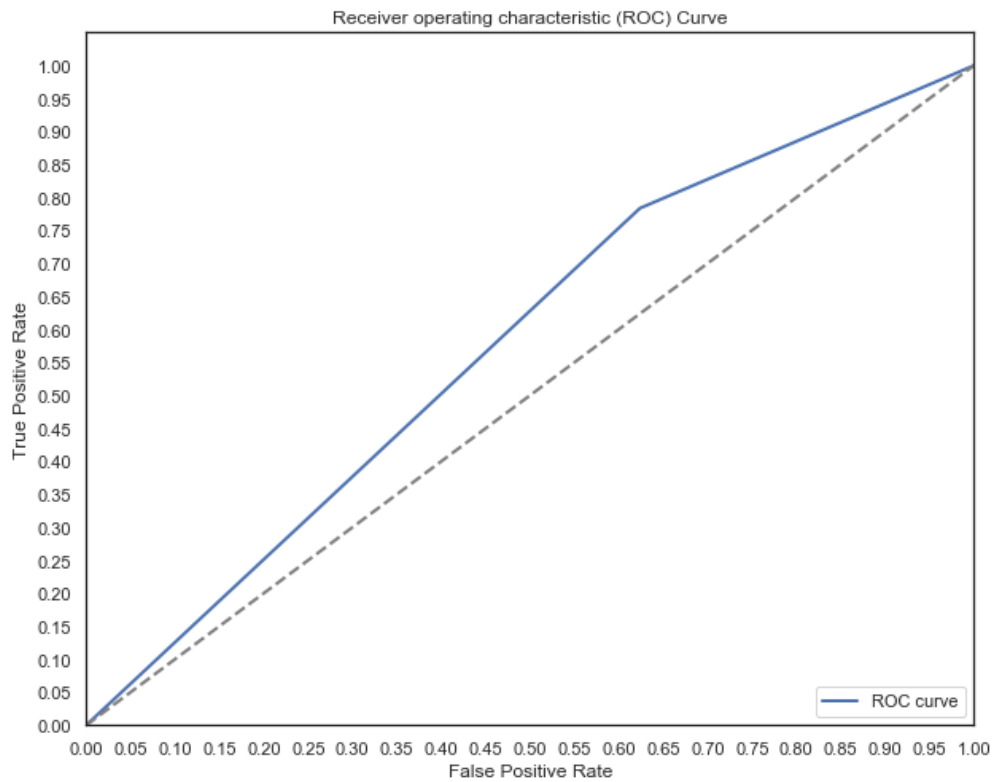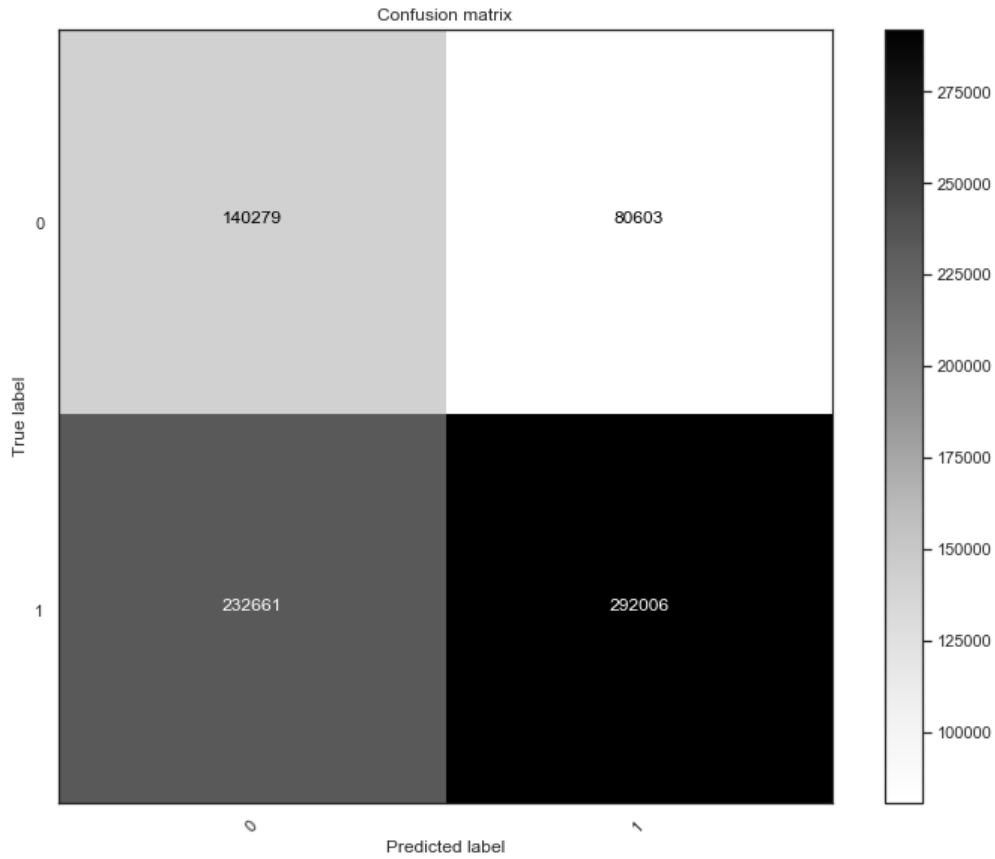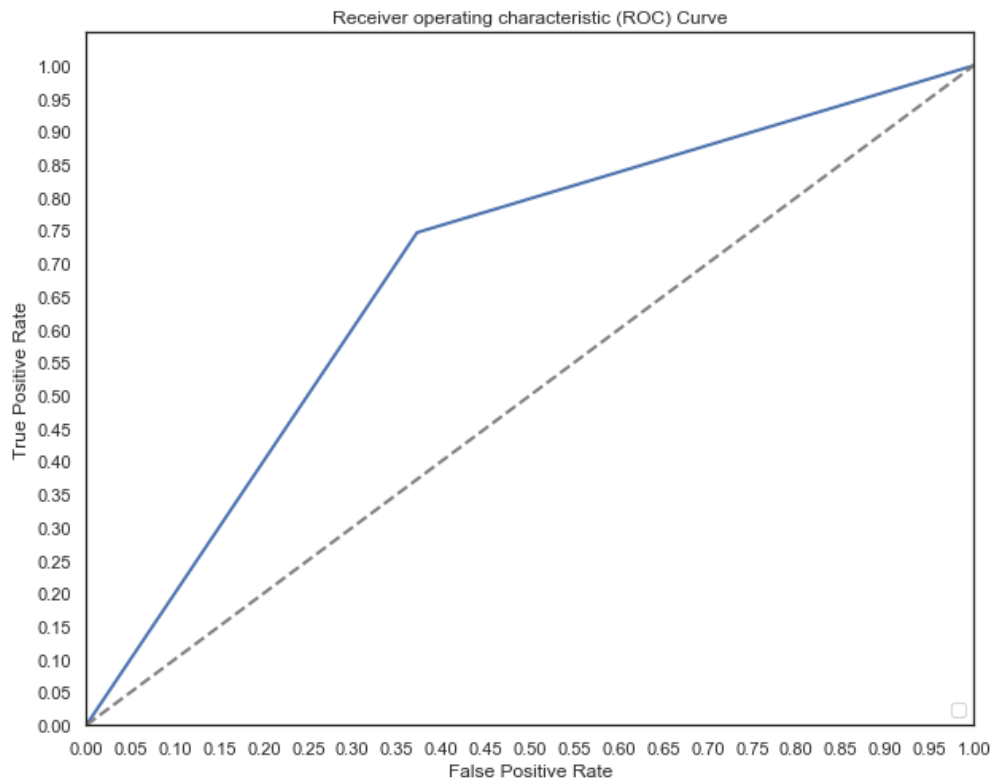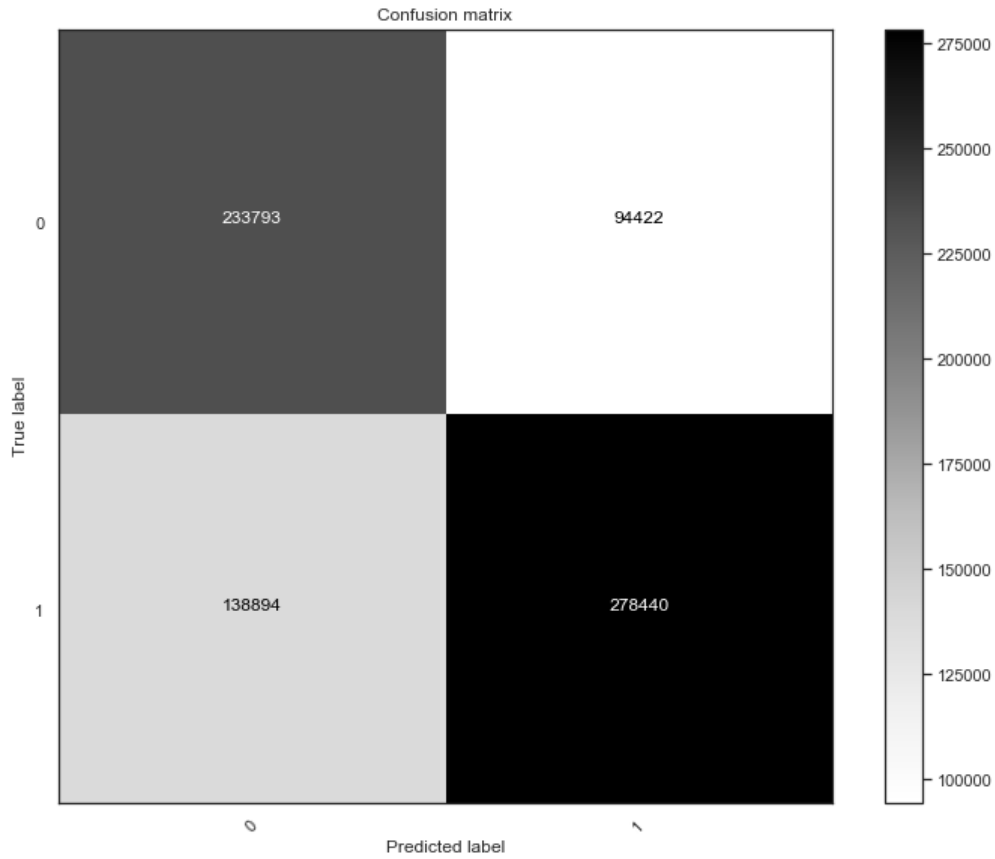
The class imbalance issue is again apparent, distorting the performance results. To address this, I'll rebalance the samples with both a SMOTE and under-sampling approach, and complete the models again.
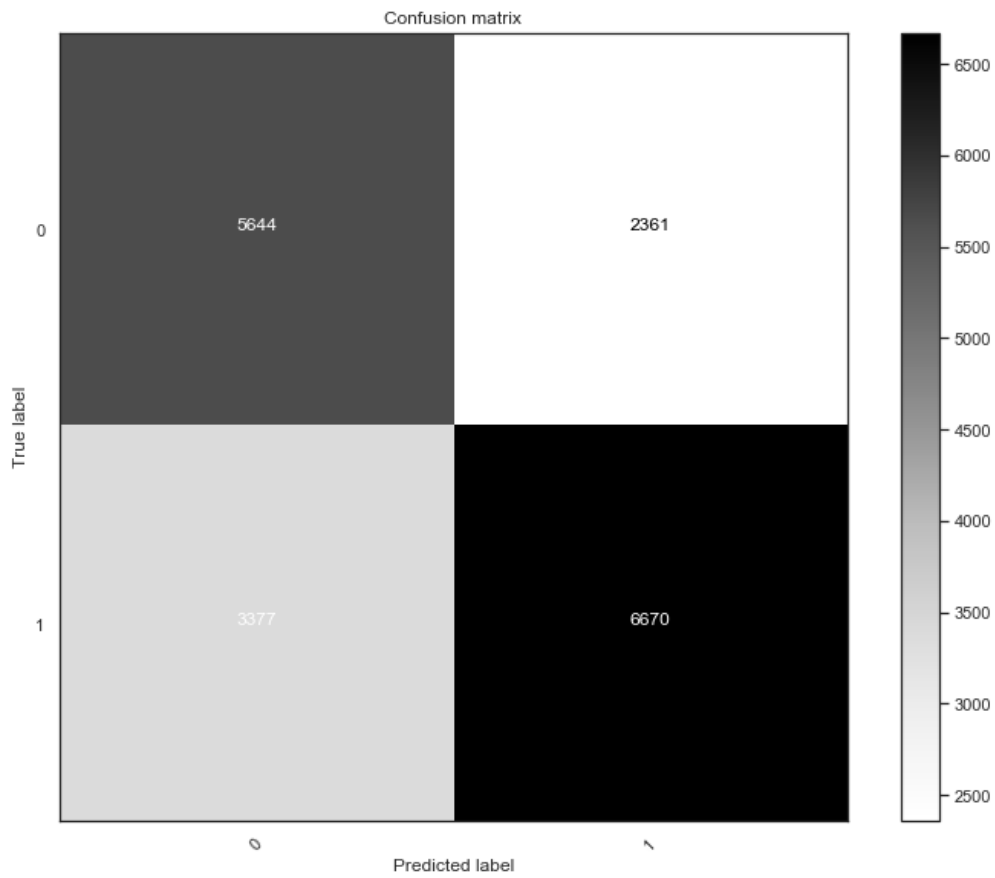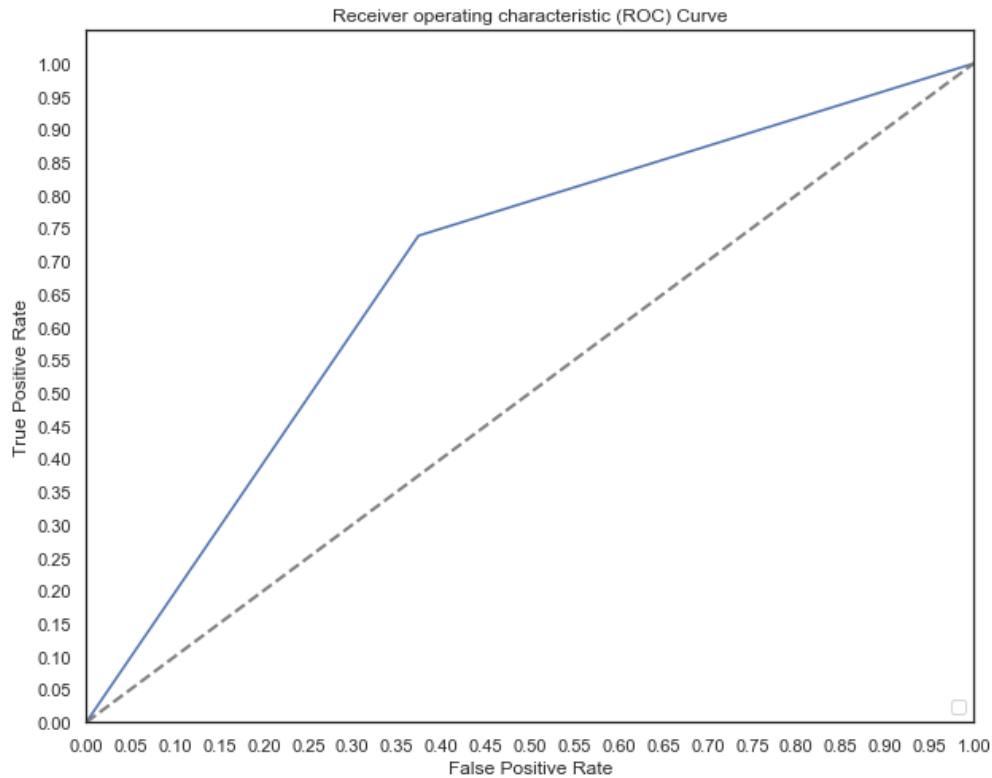
*SMOTE models with Logistic regression*

## Confusion matrix



## Receiver operating characteristic (ROC) Curve



In [ ]:

Confusion matrix

Receiver operating characteristic (ROC) Curve

*Under-sampled models with Logistic regression*

For this baseline logistic model, the udpated dataset has better performance. I'll move forward with these complete set of features as I modify the model moving forward.
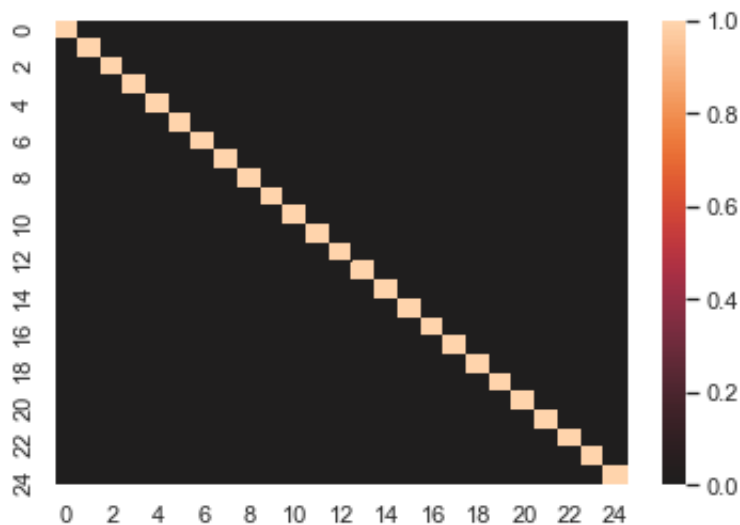


Confusion matrix

Receiver operating characteristic (ROC) Curve

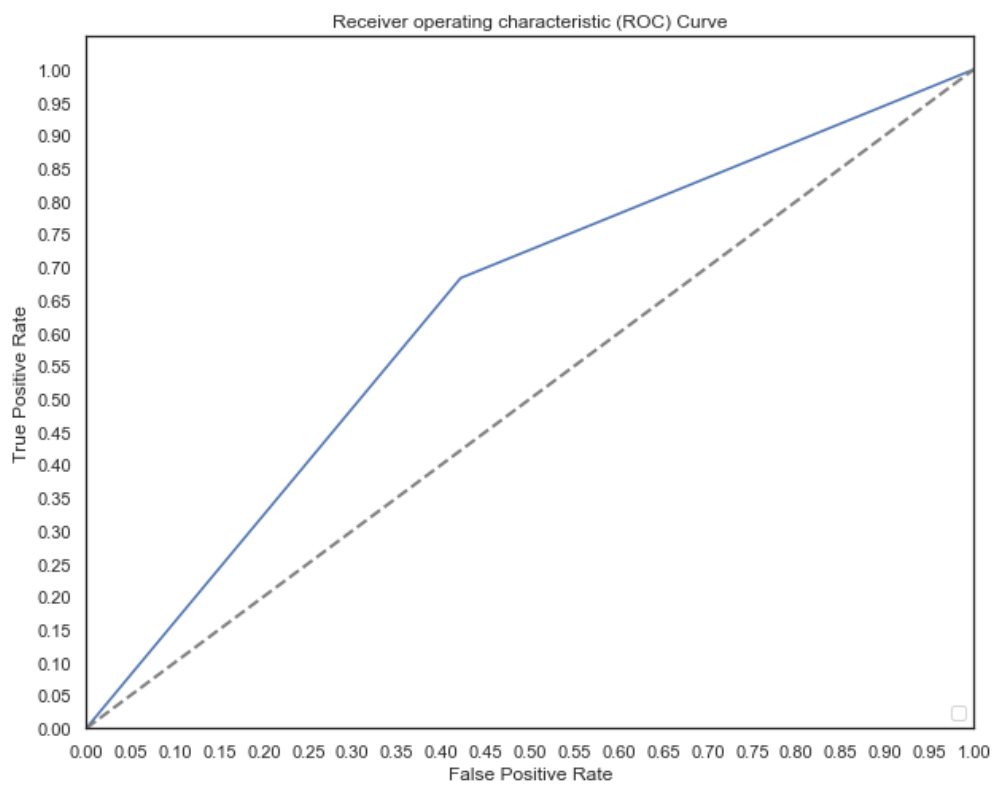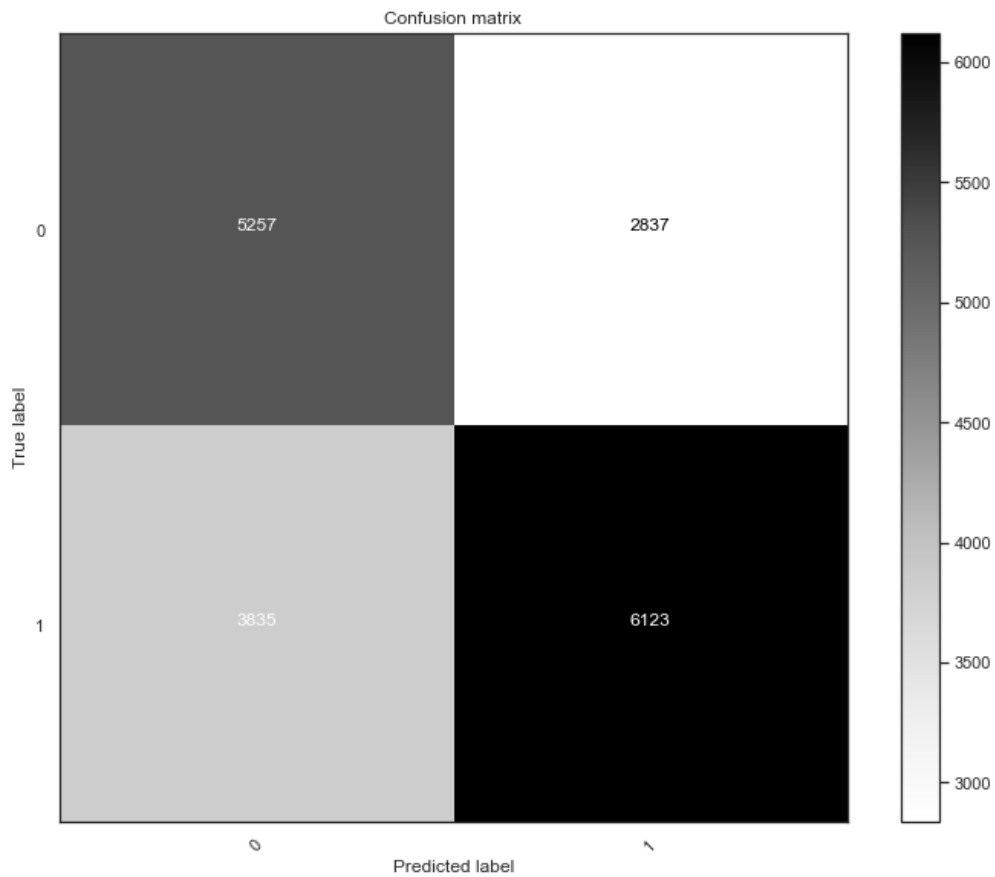*Pipeline classification models*
*Principle Component Analysis*

For this final model, to start, I implemented pricinple component analysis to reduce the number of dimensions as a precaution against issues arising from too many features. I kept 25 features in my model since it still enabled a model that could explain over 90 percent of the varioation. As you can see below, the colinearity among the remaining features is insignificant.
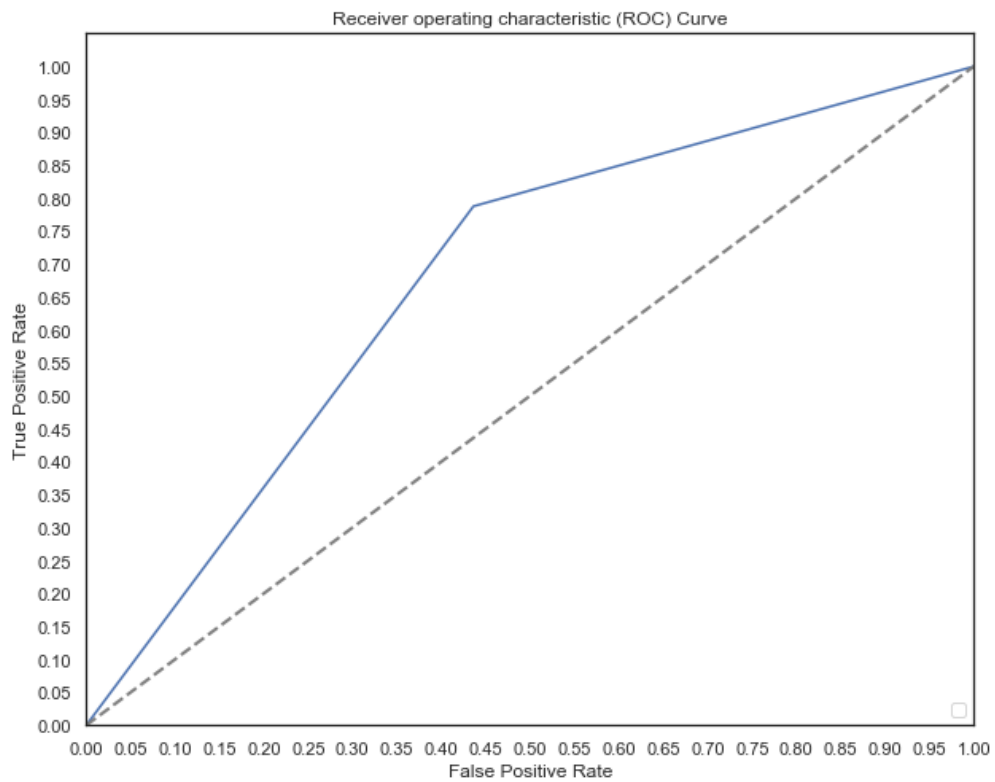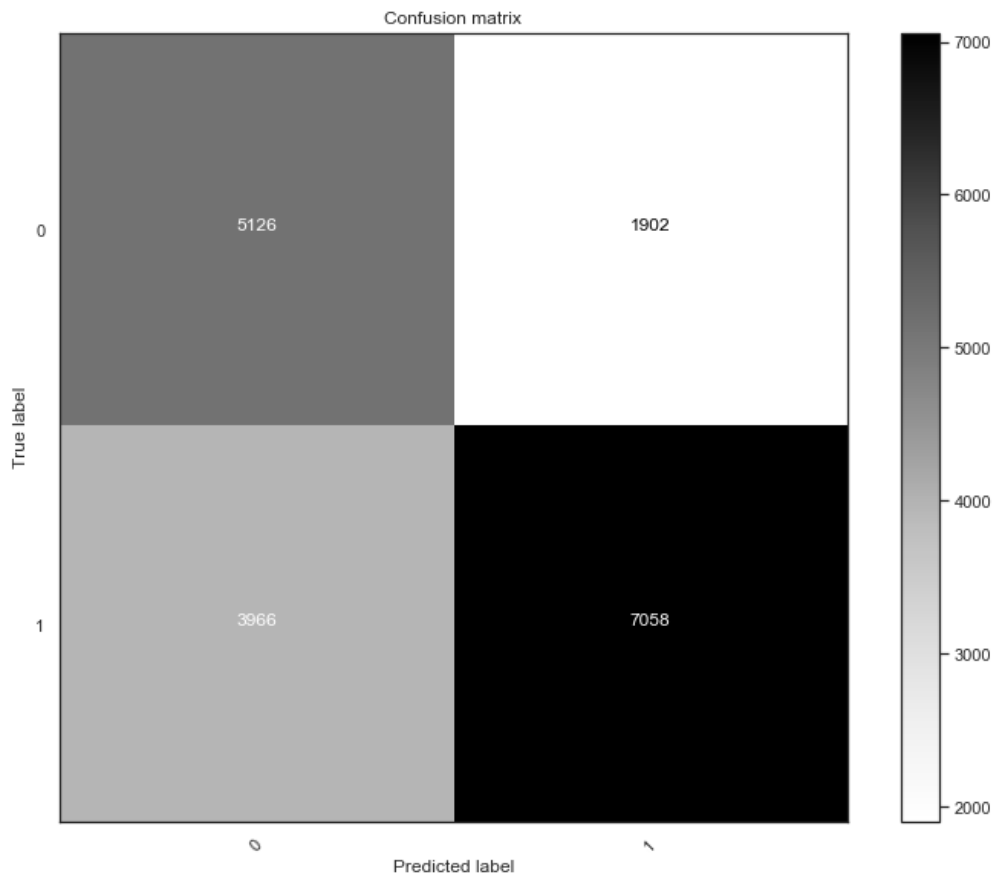


*Pipelines*

In [ ]:
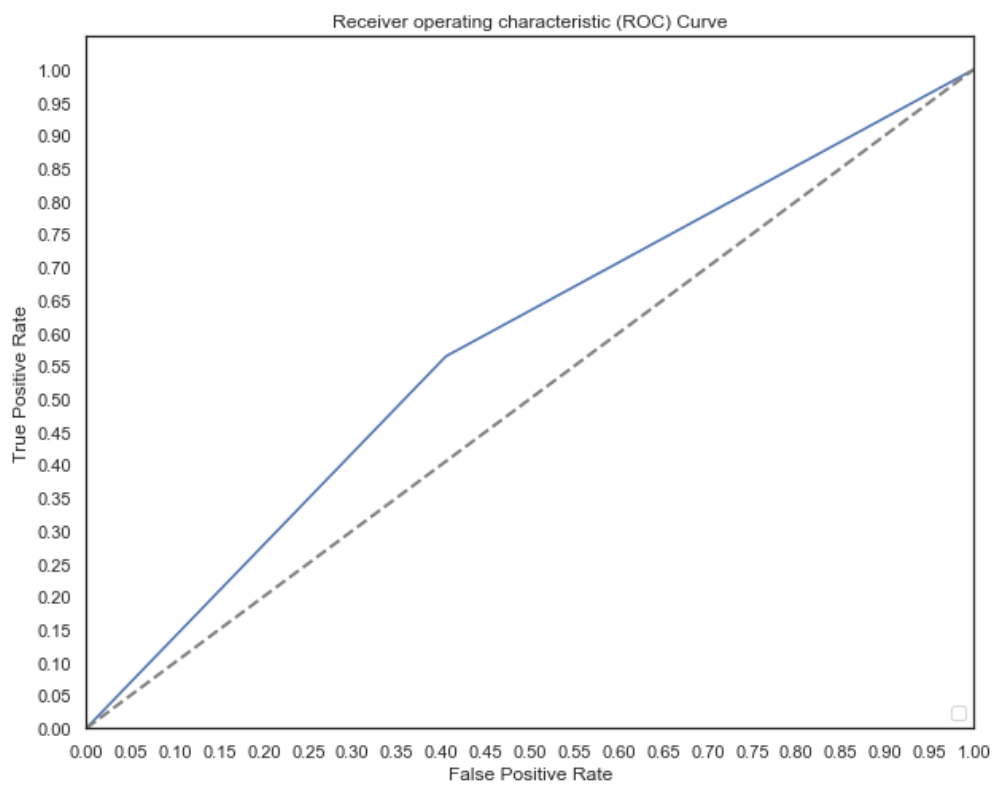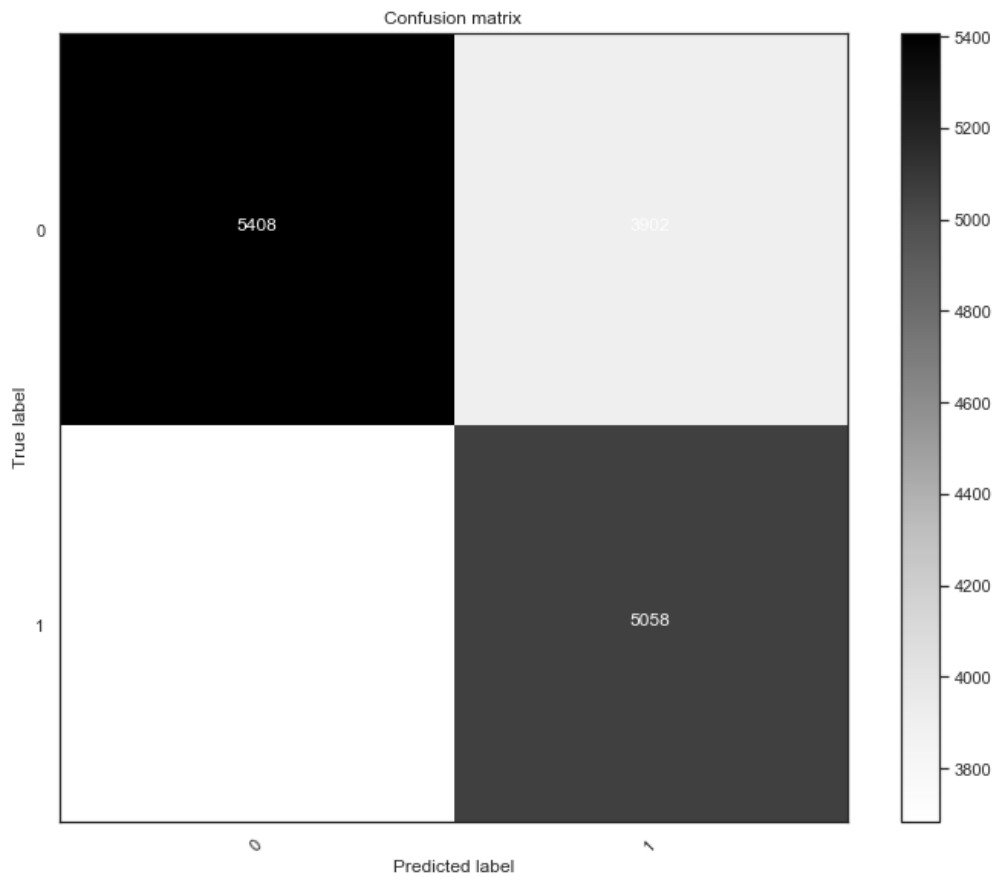
## Confusion matrix



## Receiver operating characteristic (ROC) Curve



In [ ]:

## Confusion matrix



## Receiver operating characteristic (ROC) Curve



In [ ]:

Confusion matrix



Receiver operating characteristic (ROC) Curve

In [ ]:

*Gridsearch*

Given my decision to use an udnersmpampled model, I next used grid search to find the right parameters to use for the SVM model. The gridsearch results indicated that the best parameters were:

In [ ]:

```
Best accuracy: 0.681

Best params:
 {'clf__C': 10, 'clf__gamma': 0.01, 'clf__kernel': 'rbf'}
```

In [ ]:

```
              precision    recall  f1-score   support

           0       0.74      0.57      0.64      9092
           1       0.64      0.79      0.71      8960

    accuracy                           0.68     18052
   macro avg       0.69      0.68      0.68     18052
weighted avg       0.69      0.68      0.68     18052
```



Confusion matrix