



# Predicting Unsatisfied Guests from TripAdvisor Reviews

**Author:** Albane Colmenares

**Date:** January 4th, 2024

---

## Table of Content

1. [Overview](#)
2. [Business Understanding](#)
3. [Data Understanding](#)
4. [Data Preparation](#)
5. [Modeling](#)
6. [Evaluation](#)
7. [Findings & Recommendations](#)
8. [Limits & Next Steps](#)

## 1. Overview

This notebook examines reviews about hotels from guests who stayed there and predicts whether the sentiment of reviews is positive or negative.

The organization of this notebook follows the CRoss Industry Standard Process for Data Mining (CRISP-DM) is a process model that serves as the base for a data science process.

## 2. Business Understanding

The increasing significance of online reputation and social media in the hospitality industry has become integral to a hotel's operational and financial success. The [report](https://ecornell-impact.cornell.edu/the-impact-of-social-media-on-lodging-performance/) (<https://ecornell-impact.cornell.edu/the-impact-of-social-media-on-lodging-performance/>) from Cornell's Center for Hospitality Research highlights the rising trend of consumers relying on online reviews during hotel bookings. The study reveals the link between a hotel's online reputation and its pricing, occupancy rates and revenue performance.

Recognizing the challenges faced by hotels in managing online reputation, a **new branch of TripAdvisor** is now dedicated to supporting and advising hotels on how to **improve their online reputation**. The branch launched a new B2B product identifying **unsatisfied guests in real time** with the goal to improve their experience before they leave.

Our first client is a worldwide hotel chain: Hyatt. The intended users are corporate directors of guest experience.

Through real-time guest reviews and a sophisticated algorithm, the platform alerts general managers about guests requiring immediate attention. The goal is to enable proactive measures to enhance satisfaction. This innovative tool, backed by an analysis of over 20,000 real guest reviews, also offers tailored recommendations for improvement upon implementation, aligning with the goal of turning [detractors](https://www.datasciencecentral.com/what-is-a-good-net-promoter-score-for-the-hotel-resort-industry/) (<https://www.datasciencecentral.com/what-is-a-good-net-promoter-score-for-the-hotel-resort-industry/>) into promoters.

The data is hosted on [Kaggle](https://www.kaggle.com/datasets/andrewmvrd/trip-advisor-hotel-reviews) (<https://www.kaggle.com/datasets/andrewmvrd/trip-advisor-hotel-reviews>) and is provided by **LARXEL**.

### Objectives

In order to access the file, access the data source through the link provided and click on the download button to download the dataset. This will download a zipped folder. The below code unzips it and loads it into memory.

The dataset contains about twenty thousand reviews collected from a hotel's TripAdvisor page. Each review is assigned a rating from 1 to 5, 1 being the lowest (negative review), 5 being the highest (positive review).

**Our first client** is a worldwide hotel chain, Hyatt. The intended users are corporate directors of guest experience.

Through real-time guest reviews and a sophisticated algorithm, the platform alerts general managers about guests requiring immediate attention. The goal is to enable proactive measures to enhance satisfaction. This innovative tool, backed by an analysis of over 20,000 real guest reviews, also offers tailored recommendations for improvement upon implementation, aligning with the goal of turning [detectors \(<https://www.datasciencecentral.com/what-is-a-good-net-promoter-score-for-the-hotel-resort-industry/>\)](https://www.datasciencecentral.com/what-is-a-good-net-promoter-score-for-the-hotel-resort-industry/) into promoters.

The data is hosted on [Kaggle \(<https://www.kaggle.com/datasets/andrewmvd/trip-advisor-hotel-reviews>\)](https://www.kaggle.com/datasets/andrewmvd/trip-advisor-hotel-reviews), and is provided by LARXEL.

#### Objectives

In order to access the file, access the data source through the link provided and click on the download button to save the zip file. This will download a zipped folder. The below code unzips it then loads the customized recommendations targeting areas that historically caused guest complaints.

The dataset contains about twenty thousand reviews collected from a hotel's TripAdvisor page. Each review is assigned a rating from 1 to 5, with 1 being the lowest (negative review), 5 being the highest (positive review).

- Features

Prior to preprocessing, the columns are:

- Review : the actual review's record
- Rating : a number from 1 to 5, given by the guest at the time of the review

- Target

The target will be identifying negative reviews. It will be created based on the Rating column, and will be named Sentiment . Based on a given set of reviews, I will predict if the review's sentiment was negative. To be more precise, it will be identified as detractors . This term is defined by the metric measuring online reputation: Net Promoter Score.

More on NPS below.

- Loading the data

```
In [1]: # Importing the necessary packages
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.ticker as ticker
import seaborn as sns
import warnings
from sklearn.model_selection import train_test_split
import nltk

%matplotlib inline
```

```
In [4]: # Printing the number of rows and columns in the dataset
print(f'The dataset has {len(raw_df)} rows and {len(raw_df.columns)} columns.' )
```

```
In [2]: # The dataset has 20991 rows and 2 columns
raw_df = pd.read_csv('data/tripadvisor_hotel_reviews.csv.zip', compression='zip',
```

```
In [5]: # Inspecting the number of reviews referring to each rating
```

```
In [3]: raw_df['Rating'].value_counts() of the DataFrame
```

```
Out[5]: 5    9054
```

```
Out[3]: 4    6039
       3    2184
       Review  Rating
       2    1793
       1    1421
       Name: Rating, dtype: int64
       0    nice rooms not 4* experience hotel monaco seat...
       1    nice hotel expensive parking got good deal sta...
       2    nice room special offer diamond member hilton...
       3    ok nothing special large room monaco seat...
```

```
In [6]: # Inspecting the number of reviews referring to each rating
raw_df['Rating'].value_counts(normalize=True)
```

```
Out[6]: 4    0.441853
       5    0.294715
       3    0.106583
```

The text file is encoded using Latin-1 encoding - and is open as is. Several encodings were tried to ensure the right one matched: utf-8, utf-16, ascii for example.

```
In [4]: # Printing the number of rows and columns in the dataset
print(f'The dataset has {len(raw_df)} rows and {len(raw_df.columns)} columns.' )
```

```
In [2]: # The dataset has 20491 rows and 2 as columns
raw_df = pd.read_csv('data/tripadvisor_hotel_reviews.csv.zip', compression="zip")
```

```
In [5]: # Inspecting the number of reviews referring to each rating
```

```
In [3]: raw_df['Rating'].value_counts() of the DataFrame
```

```
Out[5]: 5    9054
Out[3]: 4    6039
         3    2184
         2   1793
         1    1429
Name: Rating, dtype: int64
```

	Review	Rating
2	1793	3
1	1429	4
Name: Rating, dtype: int64	2	

```
In [6]: # Inspecting the number of reviews referring to each rating
raw_df['Rating'].value_counts(normalize=True)
```

```
Out[6]: 5    0.441853
        4    0.294715
        3    0.106583
        2    0.087502
        1    0.069348
Name: Rating, dtype: float64
```

There is a clear majority of reviews with a rating of 5: over 44% of them while all other ratings are inequally distributed. The dataset is highly imbalanced. I will review whether grouping them as a binary classification can address the issue.

## 4. Data Preparation

### 4: 1- Data Cleaning

For a better readability of the reviews, the column width will be increased.

This [link](https://www.tripadvisor.com>ShowTopic-g1-i12105-k11476502-What_is_the_maximum_character_limit_on_a_review-Tripadvisor_Support.html) ([https://www.tripadvisor.com>ShowTopic-g1-i12105-k11476502-What\\_is\\_the\\_maximum\\_character\\_limit\\_on\\_a\\_review-Tripadvisor\\_Support.html](https://www.tripadvisor.com>ShowTopic-g1-i12105-k11476502-What_is_the_maximum_character_limit_on_a_review-Tripadvisor_Support.html)) from TripAdvisor suggests the maximum number of characters for reviews is 20,000. I will set the limit to 20,000 to ensure all words are encountered for.

```
In [7]: # Increasing column width
pd.set_option('max_colwidth', 20000)
```

```
In [8]: # Inspecting the first 5 rows of the DataFrame now that there no Longer is a Limit
raw_df.head()
```

1	star ranged inamerent nor neiprui, asked desk good breakast spots neignborhood nood told no hotels, gee best breakfast spots seattle 1/2 block away convenient hotel does not know exist, arrived late night 11 pm inside run bellman busy chatng cell phone help bags.prior arrival emailed hotel inform 20th anniversary half really picky wanted make sure good, got nice email saying like deliver bottle champagne chocolate covered strawberries room arrival celebrate, told needed foam pillows, arrival no champagne strawberries no foam pillows great room view alley needed foam pillows, good by housekeeping staff cleaner room property, impressed left morning shopping room got short trips 2 hours, beds comfortable.not good ac-heat control 4 x 4 inch screen bring green shine directly eyes night, light sensitive tape controls.this not 4 start hotel clean business hotel super high rates, better chain hotels seattle	2
---	---	---

**4: 1- a) Copying the raw data**

In order to keep the raw data as it is in case it needs to be accessed in the future, an exact copy of the dataset will be made. The new DataFrame uses the same columns but gives large bathroom mediterranean suite comfortable bed pillowsattentive housekeeping staffnegatives ac unit malfunctioned stay desk disorganized, missed 3 separate wakeup calls, concierge busy hard

- **Review** touch, did n't provide guidance special requests.tv hard use ipod sound dock suite non
- **Rating** 5, decided book mediterranean suite 3 night weekend stay 1st choice rest party filled, comparison w spent 45 night larger square footage room great soaking tub whirlpool jets nice shower.before stay hotel arrange car service price 53 tip reasonable driver waiting arrival.checkin easy downside room picked 2 person jacuzzi tub no bath accessories salts bubble bath did n't

```
In [9]: # Making a copy of the raw DataFrame to check for duplicates
df = raw_df.copy()
df = df.drop_duplicates()
```

```
In [10]: # Verifying the changes appear
df.head()
```

Despite having removed some reviews, the count was verified inside the raw data, which confirms the previous were not provided entirely

saying like deliver bottle champagne chocolate covered strawberries room arrival celebrate, told needed foam pillows, arrival no champagne strawberries no foam pillows great room view alley

**4: 1- a) Copying the raw data**  
high rise building good not better housekeeping staff cleaner room property, impressed left morning shopping room got short trips 2 hours, beds comfortable.not good ac-heat control 4 x 4 inch screen bring green shine directly eyes night, light sensitive tape controls.this not 4 start hotel clean business hotel super high rates, better chain hotels seattle,

In order to keep the raw data as it is in case it needs to be accessed in the future, an exact copy of the dataset will be made. The new DataFrame has the same columns but gives large bathroom mediterranean suite comfortable bed pillowsattentive housekeeping staffnegatives ac unit malfunctioned stay desk disorganized, missed 3 separate wakeup calls, concierge busy hard

- **Review**: touch, did n't provide guidance special requests.tv hard use ipod sound dock suite non functioning.
- **Rating**: decided book mediterranean suite 3 night weekend stay 1st choice rest party filled, comparison w spent 45 night larger square footage room great soaking tub whirlpool jets nice shower.before stay hotel arrange car service price 53 tip reasonable driver waiting arrival.checkin easy downside room picked 2 person jacuzzi tub no bath accessories salts bubble bath did n't

In [9]: #<sup>2</sup> Making a copy of the raw data  
df = df.copy()  
df = df[~df.duplicated()]  
df = df[~df.duplicated(keep='last')]

In [10]: # Verifying the changes applies  
df.head()

Despite increasing the column width, some reviews are not finished. This was verified inside the raw data, while confirming the reviews were not repeated entirely.

raw data, while confirming the reviews were not repeated entirely.  
e mailed hotel inform 20th anniversary half really picky wanted make sure good, got nice email saying like deliver bottle champagne chocolate covered strawberries room arrival celebrate, told needed foam pillows, arrival no champagne strawberries no foam pillows great room view alley high rise building good not better housekeeping staff cleaner room property, impressed left morning shopping room got short trips 2 hours, beds comfortable.not good ac-heat control 4 x 4 inch screen bring green shine directly eyes night, light sensitive tape controls.this not 4 start hotel clean business hotel super high rates, better chain hotels seattle,

nice rooms not 4\* experience hotel monaco seattle good hotel n't 4\* level.positives large bathroom mediterranean suite comfortable bed pillowsattentive housekeeping staffnegatives ac unit malfunctioned stay desk disorganized, missed 3 separate wakeup calls, concierge busy hard touch, did n't provide guidance special requests.tv hard use ipod sound dock suite non functioning. decided book mediterranean suite 3 night weekend stay 1st choice rest party filled, comparison w spent 45 night larger square footage room great soaking tub whirlpool jets nice shower.before stay hotel arrange car service price 53 tip reasonable driver waiting arrival.checkin easy downside room picked 2 person jacuzzi tub no bath accessories salts bubble bath did n't stay, night got 12/11 checked voucher bottle champagne nice gesture fish waiting room, impression room huge open space felt room big, tv far away bed chore change channel, ipod dock broken disappointing.in morning way asked desk check thermostat said 65f 74 2 degrees warm try cover face night bright blue light kept, got room night no, 1st drop desk, called maintenance came look thermostat told play settings happy digital box wo n't work, asked

wake up 10am morning did n't happen, called later from nap wake up forgot. 10am wake up

## 4: 1- b) Missing data

In the next section, the missing values are inspected.

None of the two columns have null data, no modification is necessary.

In [11]: # Looking for missing values  
df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20491 entries, 0 to 20490
Data columns (total 2 columns):
 #   Column   Non-Null Count  Dtype  
--- 
 0   Review    20491 non-null   object 
 1   Rating    20491 non-null   int64  
dtypes: int64(1), object(1)
memory usage: 320.3+ KB
```

## 4: 1- d) Categorizing the ratings into a binary classification

### 4: 1- c) Handling duplicates

- **Sentiment**

No duplicate rows were identified. No changes are necessary.

In [12]: Five sentiment categories are described, which could be grouped in two: detractors and positive.  
print(df[(df[df.duplicated()]) + f' duplicate rows were identified.'])

Let's first review this assumption that the rating of 5 corresponds to a high positive review.

In [13]: # Filtering on the first 5 reviews with a rating of 5  
df[df['Rating'] == 5].head()

	Review	Rating	Review	Rating
	unique, great stay, wonderful time hotel monaco, location excellent short stroll main downtown shopping area, pet friendly room showed no signs animal hair smells, monaco suite sleeping area			
3	big striped curtains pulled closed nice touch felt cosy, goldfish named brandi enjoyed, did n't partake free wine coffee/tea service lobby thought great feature, great staff friendly, free wireless internet hotel worked suite 2 laptops, decor lovely eclectic mix patterns color palette, animal print bathrobes			

## 4: 1- d) Categorizing the ratings into a binary classification

### 4: 1- c) Handling duplicates

- Sentiment

No duplicate rows were identified. No changes are necessary.

```
In [12]: # Five sentiment categories are described, which could be grouped in two: detractors and
# not_detractors(df[df.duplicated()]) + f' duplicate rows were identified.')

Let's first review this assumption that the rating of 5 corresponds to a high positive review.
```

0 duplicate rows were identified.

```
In [13]: # Filtering on the first 5 reviews with a rating of 5
df[df['Rating'] == 5].head()
```

Out[13]:

	Review	Rating	Review	Rating
3	unique, great stay, wonderful time hotel monaco, location excellent short stroll main downtown shopping area, pet friendly room showed no signs animal hair smells, monaco suite sleeping area big striped curtains pulled closed nice touch felt cosy, goldfish named brandi enjoyed, did n't partake free wine coffee/tea service lobby thought great feature, great staff friendly, free wireless internet hotel worked suite 2 laptops, decor lovely eclectic mix patterns color palette, animal print bathrobes feel like rock stars, nice did n't look like sterile chain hotel hotel personality excellent stay,	5		
4	great stay great stay, went seahawk game awesome, downfall view building did n't complain, room huge staff helpful, booked hotels website seahawk package, no charge parking got voucher taxi, problem taxi driver did n't want accept voucher barely spoke english, funny thing speak arabic called started making comments girlfriend cell phone buddy, took second realize just said fact speak language face priceless, ass told, said large city, told head doorman issue called cab company promptly answered did n't, apologized offered pay taxi, bucks 2 miles stadium, game plan taxi return going humpin, great walk did n't mind, right christmas wonderful lights, homeless stowed away building entrances leave, police presence not greatest area stadium, activities 7 blocks pike street waterfront great coffee shops way, hotel maintained foyer awesome, wine tasting available evening, best dog, taking st. bernard time family, safes hotel located service desk room, bathroom huge jetted tub huge, funny house keeping walked girlfriend getting dressed, did n't hear knock doing turn service, screamed girlfriend screams hit floor laughing, started talking spanish worked, place recommend price, check online deals just good not better, besite contains deals vouchers travel websites n't tell,	5		
5	love monaco staff husband stayed hotel crazy weekend attending memorial service best friend husband celebrating 12th wedding anniversary, talk mixed emotions, booked suite hotel monte carlos, loaned beautiful tan-goldfish named joliet weekend visited dogs worked desk human companions, room decorated nicely couch used pillows, l'occitane bath amenities welcome sight, room quiet peaceful, wireless internet access wonderful server went morning leaving problems printing boarding passes, afternoon reception serves oenophile-satisfying wine australia scrumptious cookies, restaurant closed renovation stay finally ate food good drinks better, word caution restaurant larger person not sit booths wo n't fit, 5'6 125 lbs husband 5'9 175, table smack-against stomach couple inches space mighty uncomfortable patron larger pregnant, bad design opinion place decorated funky welcoming way metal wood handblown glass light fixtures expect seattle capital glass art industry, definitely stay reason,	5		
6	cozy stay rainy city, husband spent 7 nights monaco early january 2008. business trip chance come ride.we booked monte carlo suite proved comfortable longish stay, room 905 located street building, street noise not problem view interesting rooms building look dark alley midsection large office building, suite comfortable plenty room spread, bathroom attractive squeaky clean small comparison generous proportions sitting sleeping areas, lots comfortable seating options good lighting plenty storage clothing luggage, hotel staff friendly efficient, housekeeping staff did great job pleasant, requests responded quickly.the location quite good, easy walk pike street market seattle art museum notch shopping dining options.a positive experience,	5		
8	hotel stayed hotel monaco cruise, rooms generous decorated uniquely, hotel remodeled pacific bell building charm sturdiness, everytime walked bell men felt like coming home, secure, great single travelers, location fabulous, walk things pike market space little grocery/drug store block away, today green, bravo, 1 double bed room room bed couch separated curtain, snoring mom slept curtain, great food nearby,	5		

```
In [15]: # Filtering on the first 5 reviews with a rating of 5
At first glance, the reviews indicate how great the stay was, how comfortable the rooms were and
df[df['Rating'] == 1].head()

provide compliments about the building, decor and neighborhood. This confirms 5 star reviews are
```

positive.

	Review	Rating
15	I will inspect the 1 star review to confirm that these are negative before finally creating the sentiment column. horrible customer service hotel stay february 3rd 4th 2007 my friend picked hotel monaco appearing website online package included champagne late checkout 3 free valet gift spa weekend, friend checked room hours earlier came later, pulled valet young man just stood, asked valet open said, pull bags didn't offer help, got garment bag suitcase came car key room number says not valet, car park car street pull, left key working asked valet park car gets, went room fine bottle champagne oil lotion gift spa, dressed went came got bed noticed blood drops pillows sheets pillows, disgusted just unbelievable, called desk sent somebody 20 minutes later, swapped sheets left apologizing, sunday morning called desk speak management sheets aggravated rude, apparently no manager kind supervisor weekend wait monday morning, young man spoke said cover food adding person changed sheets said fresh blood rude tone, checkout 3pm package booked, 12 1:30 staff maids tried walk room opening door apologizing closing, people called saying check 12 remind package, finally packed things went downstairs check, quickly signed paper took, way took closer look room, unfortunately covered food offered charged valet, called desk ask charges lady answered snapped saying aware problem experienced monday like told earlier, life treated like hotel, not sure hotel constantly problems lucky ones stay recommend anybody know,	1
	noise airconditioner-a standard, arranged stay travel agency unfortunately warwick seattle hotel disappointment trip, 3 night stay warwick changed 3 rooms, starting minute stay hotel personnel didn't make feel like guest like intruder reluctant help solve complaints hotel right downtown 5	

In [15]: # Filtering on the first 5 reviews with a rating of 5  
At first glance, the reviews indicate how great the stay was, how comfortable the rooms were and provide compliments about the building, decor and neighborhood. This confirms 5 star reviews are

Out[15]: positive.

		Review	Rating
1		I will inspect the 1 star review to confirm that these are negative before finally creating the	
15		horrible customer service hotel stay february 3rd 4th 2007 my friend picked hotel monaco	
32		sentiment column	

weekend, friend checked room hours earlier came later, pulled valet young man just stood, asked valet open said, pull bags didn't offer help, got garment bag suitcase came car key room number says not valet, car park car street pull, left key working asked valet park car gets, went room fine bottle champagne oil lotion gift spa, dressed went came got bed noticed blood drops pillows sheets pillows, disgusted just unbelievable, called desk sent somebody 20 minutes later, swapped sheets left apologizing, sunday morning called desk speak management sheets aggravated rude, apparently no manager kind supervisor weekend wait monday morning, young man spoke said cover food adding person changed sheets said fresh blood rude tone, checkout 3pm package booked, 12 1:30 staff maids tried walk room opening door apologizing closing, people called saying check 12 remind package, finally packed things went downstairs check, quickly signed paper took, way took closer look room, unfortunately covered food offered charged valet, called desk ask charges lady answered snapped saying aware problem experienced monday like told earlier, life treated like hotel, not sure hotel constantly problems lucky ones stay recommend anybody know,

noise airconditioner-a standard, arranged stay travel agency unfortunately warwick seattle hotel disappointment trip, 3 night stay warwick changed 3 rooms, starting minute stay hotel personnel didn't make feel like guest like intruder, reluctant help solve complaints, hotel right downtown 5 minutes really good restaurants like good thing, availability room offered 2nd floor window

Indeed, the first word of the first 1 star review indicate how *horrible* the customer service was. The ranking from 1 to 5 is confirmed to be from low to high.

#### • NPS Score

The Net Promoter Score (NPS) is a metric used to measure customer satisfaction and loyalty with a hotel, (or more generally a product, service, or company). It is based on the question: "How likely is it that you would recommend our hotel? The responses can be measured on a scale from 0 to 10 or in our case: 1 to 5.



(<https://textexpander.com/blog/how-to-calculate-nps>)

According to the 5 point [NPS Scale \(<https://www.surveysensum.com/blog/11-point-and-5-point-nps-scale#:~:text=5%2Dpoint%20NPS%20Scale,-Similar%20to%20the&text=The%20respondents%20on%20the%205,why%20it%20is%20not%20rec>](https://www.surveysensum.com/blog/11-point-and-5-point-nps-scale#:~:text=5%2Dpoint%20NPS%20Scale,-Similar%20to%20the&text=The%20respondents%20on%20the%205,why%20it%20is%20not%20rec)

1. Promoters rate 5 star reviews
2. Passive customers rate 4 stars
3. Detractors rate from 1 to 3 star reviews

In [16]: # Reminding the current percentage for each star reviews given  
This indicates that the customers who rank a hotel from 1 to 3 star reviews may really damage a hotel's online reputation and negatively influence other guests in their willingness to book your hotel.  
# Storing values for promoters and detractors  
promoters = df['Rating'].value\_counts(normalize=True).get(5, 0)  
detractors = df['Rating'].value\_counts(normalize=True).get(1, 0) + df['Rating'].get(2, 0) + df['Rating'].get(3, 0)  
nps\_score = promoters / detractors  
minus the percentage of detractors (guests who do not recommend the hotel 1 to 3 star reviews)

In [17]: # Calculating the current NPS score  
As a consequence, I will categorize the guests who gave 1 to 3 star reviews as detractors and the others as not detractors {nps\_score: .2%})  
Your NPS Score is 17.84%

#### • NPS Calculation

If a hotel was to decrease the number of detractors by only 10%, the new NPS score for your hotel chain would be:

As a starting point the NPS for this dataset is calculated.

In [18]: # Calculating the potential percentage of detractors

```
In [16]: # Reminding the current percentage for each star reviews given
This indicates that the customers who rank a hotel from 1 to 3 star reviews may really damage a
df['Rating'].value_counts(normalize=True)
hotel's online reputation and negatively influence other guests in their willingness to book your
hotel
# Storing values for promoters and detractors
promoters = df['Rating'].value_counts(normalize=True).get(5, 0)
NPS is calculated as df['Rating'].value_counts(normalize=True).get(1, 0) + df['Rating'].value_counts(normalize=True).get(2, 0)
• the percentage of promoters (guests who highly recommend the hotel: 5 star reviews)
• the percentage of detractors (guests who do not recommend the hotel: 1 to 3 star
reviews)
```

```
In [17]: # Calculating the current NPS score
As a consequence I will categorize the guests who gave 1 to 3 star reviews as detractors and
the others as not detractors {nps_score:.2%})
print(f'Your NPS Score is {nps_score:.2%}')
```

Your NPS Score is 17.84%

- **NPS Calculation**

If a hotel was to decrease the number of detractors by only 10%, the new NPS score for your hotel chain would be:

As a starting point the NPS for this dataset is calculated.

```
In [18]: # Calculating the potential percentage of detractors
potential_detractors = detractors * 0.9
```

```
In [19]: # Calculating the potential NPS Score
potential_nps_score = promoters - potential_detractors
print(f'Your potential NPS Score could be {potential_nps_score:.2%}')
```

Your potential NPS Score could be 20.48%

```
In [20]: # Calculating the growth in NPS Score after reducing detractors by 10%
growth_in_nps = (potential_nps_score - nps_score) / nps_score
print(f'By decreasing only the percentage of detractors in your hotels by only 10% you
could increase your NPS Score by {growth_in_nps:.2%}!')
```

By decreasing only the percentage of detractors in your hotels by only 10%, you could increase your NPS Score by 14.76%!

I will create this new category under the column Sentiment .

```
In [21]: # Creating the column sentiment
df['Sentiment'] = df['Rating'].apply(lambda x: 'detractors' if x in [1, 2, 3] else 'not_detractors')
```

```
In [22]: # Verifying the column was created
df.head()
```

		Review	Rating	Sentiment
0	nice hotel expensive parking got good deal stay hotel anniversary, arrived late evening took advice previous reviews did valet parking, check quick easy, little disappointed non-existent view room room clean nice size, bed comfortable woke stiff neck high pillows, not soundproof like heard music room night morning loud bangs doors opening closing hear people talking hallway, maybe just noisy neighbors, aveda bath products nice, did not goldfish stay nice touch taken advantage staying longer, location great walking distance shopping, overall nice experience having pay 40 parking night,			4 not_detractors
1	ok nothing special charge diamond member hilton decided chain shot 20th anniversary seattle, start booked suite paid extra website description not, suite bedroom bathroom standard hotel room, took printed reservation desk showed said things like tv couch ect desk clerk told oh mixed suites description kimpton website sorry free breakfast, got kidding, embassy suits sitting room bathroom bedroom unlike kimpton calls suite, 5 day stay offer correct false advertising, send kimpton preferred guest website email asking failure provide suite advertised website reservation description furnished hard copy reservation printout website desk manager duty did not reply solution, send email trip guest survey did not follow email mail, guess tell concerned guest the staff ranged indifferent not helpful, asked deck need breakfast note neighborhood hand told no hotels see heat breakfast			2 detractors

In [22]: # Verifying the column was created  
df.head()

Out[22]:

		Review	Rating	Sentiment
0	nice hotel expensive parking got good deal stay hotel anniversary, arrived late evening took advice previous reviews did valet parking, check quick easy, little disappointed non-existent view room room clean nice size, bed comfortable woke stiff neck high pillows, not soundproof like heard music room night morning loud bangs doors opening closing hear people talking hallway, maybe just noisy neighbors, aveda bath products nice, did not goldfish stay nice touch taken advantage staying longer, location great walking distance shopping, overall nice experience having pay 40 parking night,		4	not_detractors
1	ok nothing special charge diamond member hilton decided chain shot 20th anniversary seattle, start booked suite paid extra website description not, suite bedroom bathroom standard hotel room, took printed reservation desk showed said things like tv couch ect desk clerk told oh mixed suites description kimpton website sorry free breakfast, got kidding, embassy suits sitting room bathroom bedroom unlike kimpton calls suite, 5 day stay offer correct false advertising, send kimpton preferred guest website email asking failure provide suite advertised website reservation description furnished hard copy reservation printout website desk manager duty did not reply solution, send email trip guest survey did not follow email mail, guess tell concerned guest, the staff ranged indifferent not helpful, asked desk good breakfast spots neighborhood hood told no hotels, gee best breakfast spots seattle 1/2 block away convenient hotel does not know exist, arrived late night 11 pm inside run bellman busy chatting cell phone help bags prior arrival emailed hotel inform 20th anniversary half really picky wanted make sure good, got nice email saying like deliver bottle champagne chocolate covered strawberries room arrival celebrate, told needed foam pillows, arrival no champagne strawberries no foam pillows great room view alley high rise building good not better housekeeping staff cleaner room property, impressed left morning shopping room got short trips 2 hours, beds comfortable, not good ac-heat control 4 x 4 inch screen bring green shine directly eyes night, light sensitive tape controls, this not 4 start hotel clean business hotel super high rates, better chain hotels seattle,		2	detractors
2	nice rooms not 4* experience hotel monaco seattle good hotel n't 4* level, positives large bathroom mediterranean suite comfortable bed pillows attentive housekeeping staff negatives ac unit malfunctioned stay desk disorganized, missed 3 separate wakeup calls, concierge busy hard touch, did n't provide guidance special requests, tv hard use ipod sound dock suite non functioning, decided book mediterranean suite 3 night weekend stay 1st choice rest party filled, comparison w spent 45 night larger square footage room great soaking tub whirlpool jets nice shower, before stay hotel arrange car service price 53 tip reasonable driver waiting arrival, checkin easy downside room picked 2 person jacuzzi tub no bath accessories salts bubble bath did n't stay, night got 12/1a checked voucher bottle champagne nice gesture fish waiting room, impression room huge open space felt room big, tv far away bed chore change channel, ipod dock broken disappointing, in morning way asked desk check thermostat said 65f 74 2 degrees warm try cover face night bright blue light kept, got room night no, 1st drop desk, called maintenance came look thermostat told play settings happy digital box wo n't work, asked wakeup 10am morning did n't happen, called later 6pm nap wakeup forgot, 10am wakeup morning yep forgotten, the bathroom facilities great room surprised room sold whirlpool bath tub n't bath amenities, great relax water jets going,		3	detractors
3	unique, great stay, wonderful time hotel monaco, location excellent short stroll main downtown shopping area, pet friendly room showed no signs animal hair smells, monaco suite sleeping area big striped curtains pulled closed nice touch felt cosy, goldfish named brandi enjoyed, did n't partake free wine coffee/tea service lobby thought great feature, great staff friendly, free wireless internet hotel worked suite 2 laptops, decor lovely eclectic mix patterns color palette, animal print bathrobes feel like rock stars, nice did n't look like sterile chain hotel hotel personality excellent stay,		5	not_detractors
4	great stay great stay, went seahawk game awesome, downfall view building did n't complain, room huge staff helpful, booked hotels website seahawks package, no charge parking got voucher taxi, problem taxi driver did n't want accept voucher barely spoke english, funny thing speak arabic called started making comments girlfriend cell phone buddy, took second realize just said fact speak language face friendless, as told said large city, told head doorman issue called cab company promptly answer did n't, apologized offered pay taxi, bucks 2 miles stadium, game plan taxi return going humpin great walk did n't mind night christmas wonderful lights, homes less stowed away building entrances leave, police presence not greatest area stadium, activities 7 blocks pike street waterfront great coffee shops way, hotel maintained by awesome, wine tasting available evening, best dog, detrac tors king st, bernard 11263401, safes hotel located service desk room, bathroom Name:娼妓卖淫, value_counts().count normalize=True)		5	not_detractors

In [24]: # And of course

# Inspecting the new value counts for the Sentiment category  
df['Sentiment'].value\_counts()

Out[24]: not\_detractors 15093  
detractors 5398  
Name: Sentiment, dtype: int64

The grouping also started addressing the class imbalance. Nevertheless, the dataset remains highly imbalanced.

In [23]: # Inspecting the new value counts for the Sentiment category  
df['Sentiment'].value\_counts()  
# will visually represent the split of number of reviews.

Out[23]: not\_detractors 15093  
detractors 5398  
Name: Sentiment, dtype: int64

```
In [24]: # And priceless, ass told us
# Inspecting the new value counts for the Sentiment category
df['Sentiment'].value_counts()
```

5 not\_detractors

```
Out[24]: not_detractors    15093
detractors        5398
Name: Sentiment, dtype: int64
```

The grouping also started addressing the class imbalance. Nevertheless, the dataset remains highly imbalanced.

```
In [23]: # Inspecting the new value_counts for the Sentiment category
df['Sentiment'].value_counts()
I will visually represent the split of number of reviews.
```

```
Out[23]: not_detractors    15093
detractors        5398
Name: Sentiment, dtype: int64
```

```
In [25]: from matplotlib.ticker import FuncFormatter
# Creating a function to plot the number of reviews by sentiment

# Defining custom colors
custom_colors = ['#00AF87', '#AA4255']

def plot_sentiment_count(df, x, label_rotation):
    # Defining the figure
    fig, ax = plt.subplots(figsize=(8,6))

    # Setting facecolor
    fig.set_facecolor('#2F4858')
    ax.set_facecolor('#2F4858')

    # Plotting the count of reviews
    sns.countplot(data=df, x=x, order=df[x].value_counts().index, palette=custom_
    # Defining the labels and titles
    ax.set_xlabel(xlabel = 'Sentiment', fontsize=12, color='white')
    ax.set_ylabel(ylabel = 'Number of Reviews', fontsize=12, color='white')
    ax.set_title(f'Number of Reviews per Sentiment', fontsize=15, color='white')
```

```
In [25]: from matplotlib.ticker import FuncFormatter
# Creating a function to plot the number of reviews by sentiment

# Defining custom colors
custom_colors = ['#00AF87', '#AA4255']

def plot_sentiment_count(df, x, label_rotation):
    # Defining the figure
    fig, ax = plt.subplots(figsize=(8,6))

    # Setting facecolor
    fig.set_facecolor('#2F4858')
    ax.set_facecolor('#2F4858')

    # Plotting the count of reviews
    sns.countplot(data=df, x=x, order=df[x].value_counts().index, palette=custom_colors)

    # Defining the labels and titles
    ax.set_xlabel(xlabel = 'Sentiment', fontsize=12, color='white')
    ax.set_ylabel(ylabel = 'Number of Reviews', fontsize=12, color='white')
    ax.set_title(f'Number of Reviews per Sentiment', fontsize=15, color='white')
    ax.tick_params(axis='y', colors='white')

    # Recording the sentiment values to use them as labels
    sentiment_labels = df[x].unique()

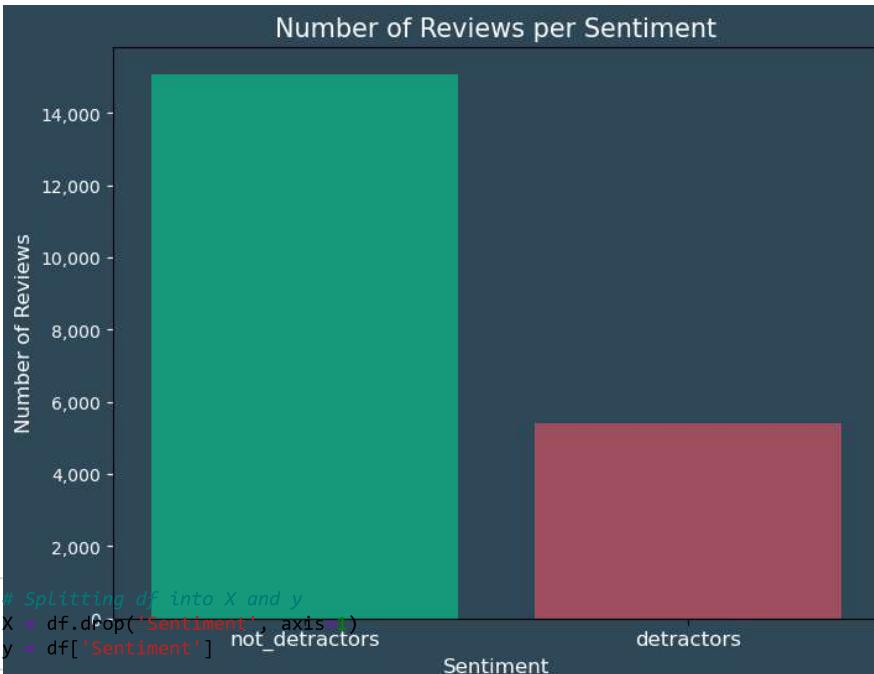
    # Setting the x-axis tick labels to the sentiment labels
    ax.set_xticks(range(len(sentiment_labels)))
    ax.set_xticklabels(labels=sentiment_labels, rotation=label_rotation, fontsize=10, color='white')

    # Formatting the y-axis labels to show thousands
    ax.yaxis.set_major_formatter(FuncFormatter(lambda x, pos: '{:,}'.format(int(x/1000)*1000)))

    # Saving the plot as a PNG with a transparent background
    plt.savefig('images/reviews_per_sentiment.png')

    # Showing the plot
    plt.show()

plot_sentiment_count(df, 'Sentiment', 0)
```



```
In [26]: # Splitting df into X and y
X = df.drop(['Sentiment'], axis=1)
y = df['Sentiment']
```

```
In [27]: # Performing a train test split
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42, stratify=y)
```

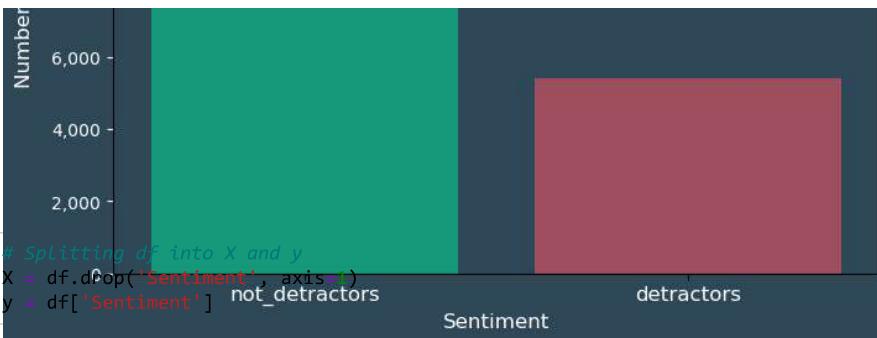
#### 4: 1-e) Performing a Train-Test Split

```
In [28]: # Inspecting the X_train data
```

The dataset is being divided into two separate subsets: a training set, and a testing (or validation) set. The validation set will allow to assess the performance of the model.

```
Out[28]:
```

Review	Rating
--------	--------



In [27]: # Performing a train test split  
X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, random\_state=42, stratify=y)

**4: 1-e) Performing a Train-Test Split**

In [28]: # Inspecting the X\_train data  
The dataset is being divided into two separate subsets: a training set, and a testing (or validation) set. The validation set will allow to assess the performance of the model.

		Review	Rating
1282	gem, pleasantly surprised accomodations helpful attentive staff need, breakfast good cookies, room lovely clean comfortable, room fronting bush street bit noisy did not away enjoyment nice boutique hotel, definitely stay time,		4
10661	loved fita wife spent nights hotel fita march 2008. hotel staff fantastic helpful pleasant, recommend hotel traveling amsterdam.the hotel located minute walk van gogh rijks museums 10 minute walk leidsplein variety stores shops restaurants located.there tram stop directly hotel, easy access entire city, did walk fita way center city 30 minutes.beautiful accesible friendly definately travel fita wonderful amsterdam,		5
17908	great modern hotel fantastic price stayed just night beginning april, alot reviews say bad area really did n't, literally 100m away nice redeveloped area 10min walk beach.me partner arrived late barcelona decided walk hotel bus station order barings, quite long walk 50mins roughly really great thing order sights tourists dont usually, arrived check quick painless room really really nice, fantastic modern design little extras want, order barcelona partner opted bus turistic, stop not far hotel goes virtually, worth short visit.i wish spent longer hotel, wish swimming pool roof open, overall fantastic hotel price paid room dont distance centre transport options,		5
11685	great experiance, stayed family golden tulip week summer, rooms somewhat small beds incredibly comfortable bathrooms good size, staff really friendly hotel clean food restaurant fantastic not bit expensive, 10 minute train ride puts amsterdam easy use, took bar nights drinks coffee great accomodating large group people, not greatest neighborhood 5 minute walk takes shops, stay,		4
9059	excellent location hotel improved, things location hotel excellent, plenty bars shops restaurants activities nearby 20 minutes walk centre amsterdam red light district interested having look, lobby area nice staff friendly checked quickly approaching room worried area hotel older described poor condition shabby, carpets dirty walls paper hanging really bad smell stale smoke drink food onions entered room group consisted married couple rooms given abysmal, dirty smelly definately not 4 star standard, hotel needs refurbishment.. sooner better, bed quarter size double just room 2 barely floor foot bed walk sides furniture old shabby badly damaged looked like selected junk shop, actually felt quite insulted given room appalling standard absolutely no way earth prepared spend night expect complained, told receptionist not room disgusting insulted fact clearly expected sleep, way apology offered executive room newer hotel discounted rate 25 euro night room apparently usually 45 euros extra room night, vast improvement expect standard, decent sized bed king matching furniture tidy bathroom nice decor plenty room, bright airy clean, looked fairly new given room place not cause complain rated 4 star hotel, newer hotel lovely pleasant bar area, staff friendly overall hard pressed better location.so note book make sure not given room old building disappointed,		3

In [29]: # Inspecting the y\_train data  
y\_train.head()

In [31]: # Inspecting the distribution in the test data  
Out[29]: 1282 not\_detractors  
10661 not\_detractors

Out[31]: 17908 detractor  
11685 not\_detractor  
Name: Sentiment, dtype: float64  
Name: Sentiment, dtype: object

## 4: 1-d) Encoding Target

In [30]: In order to make predictions useable in calculations, I will encode the target:  
# Inspecting the distribution in the train data  
y\_train.value\_counts(normalize=True)

Out[30]: not\_detractor 0.70596  
detractors 0.263404  
Name: Sentiment, dtype: float64

In [32]: # Adding in labels for filtering before making any encoding  
X\_train['label'] = [y\_train[val] for val in X\_train.index]  
X\_train[:3]

```
y_train.value_counts()
# Inspecting the distribution in the test data
In[30]: 1282    not_detractors
          10661   not_detractors
Out[31]: 17908   not_detractor 0.988483
          11320   detractor_got_detractor 0.011517
          Name: Sentiment, dtype: float64
          Name: Sentiment, dtype: object
```

## 4: 1) Encoding Target

In [30]: In order to make predictions useable in calculations, I will encode the target:  
`# Inspecting the distribution in the train data`

```
y_train.value_counts(normalize=True)
# 0 will replace not_detractors
```

```
Out[30]: not_detractor 0.736596
          detractors 0.263404
          Name: Sentiment, dtype: float64
```

In [32]: # Adding in Labels for filtering before making any encoding  
`X_train['label'] = [y_train[val] for val in X_train.index]`  
`X_train[:3]`

Out[32]:

		Review	Rating	label
1282	gem, pleasantly surprised accomodations helpful attentive staff need, breakfast good cookies, room lovely clean comfortable, room fronting bush street bit noisy did not away enjoyment nice boutique hotel, definitely stay time,		4	not_detractors
10661	loved fita wife spent nights hotel fita march 2008. hotel staff fantastic helpful pleasant, recommend hotel traveling amsterdam.the hotel located minute walk van gogh rijks museums 10 minute walk leidsplein variety stores shops restaurants located.there tram stop directly hotel, easy access entire city, did walk fita way center city 30 minutes.beautiful accesible friendly definately travel fita wonderful amsterdam,		5	not_detractors
17908	great modern hotel fantastic price stayed just night beginning april, alot reviews say bad area really did n't, literally 100m away nice redeveloped area 10min walk beach.me partner arrived late barcelona decided walk hotel bus station order barings, quite long walk 50mins roughly really great thing order sights tourists dont usually, arrived check quick painless room really really nice, fantastic modern design little extras want, order barcelona partner opted bus turistic, stop not far hotel goes virtually, worth short visit.i wish spent longer hotel, wish swimming pool roof open, overall fantastic hotel price paid room dont distance centre transport options,		5	not_detractors

In [33]: # Adding in Labels for filtering before making any encoding  
`X_test['label'] = [y_test[val] for val in X_test.index]`  
`X_test[:3]`

Out[33]:

		Review	Rating	label
5372	great location good hotel 1night stay recommend hotel want just night accommodation ahead day flight, sheraton stars mentioned architecture n't good hotel, stay night try solution night stay choose directly connected airport, pretty expensive cola 0,2l juices 4 euro price minibar pizza 12 euro ok person problems finding signs sheraton, problems big frankfurt airport.actually easy hotel airport directly connected airport terminal 1.from departure hall upstairs 2 pedestrian bridges whichconnect terminal hotel long distance trains 1 hall b near lufthansa 1 hall b c near american airlines follow signs hotel fernbahnhof long distance trains,		4	not_detractors

In [34]: # Inspecting common areas, arriving hotel shuttle available departing odd hour print('BEFORE')
print(y\_train['label'].value\_counts(normalize=True))
# Encoding sticks.
y\_train = y\_train.map(lambda x: 0 if x == 'not\_detractors' else 1)

```
# Verifying it is one hot encoded
print('After')
print(y_train['label'].value_counts(normalize=True))
```

```
Before
not_detractors    11320
detractors       4048
Name: Sentiment, dtype: int64
not_detractors    0.736596
detractors       0.263404
Name: Sentiment, dtype: float64
After
0    11320
```

```
In [34]: # Inspecting y_train
print('Before')
print(y_train.value_counts())
print(y_train.value_counts(normalize=True))

# Encoding y_train
y_train = y_train.apply(lambda x: 0 if x == 'not_detractors' else 1)

# Verifying it correctly applied
print('After')
print(y_train.value_counts())
print(y_train.value_counts(normalize=True))

Before
not_detractors    11320
detractors        4048
Name: Sentiment, dtype: int64
not_detractors    0.736596
detractors        0.263404
Name: Sentiment, dtype: float64
After
0    11320
1    4048
Name: Sentiment, dtype: int64
0    0.736596
1    0.263404
Name: Sentiment, dtype: float64
```

```
In [35]: # Inspecting y_test
print('Before')
print(y_test.value_counts())
print(y_test.value_counts(normalize=True))

# Encoding y_test
y_test = y_test.apply(lambda x: 0 if x == 'not_detractors' else 1)

# Verifying it correctly applied
print('After')
print(y_test.value_counts())
print(y_test.value_counts(normalize=True))
```

```
Before
not_detractors    3773
detractors        1350
Name: Sentiment, dtype: int64
not_detractors    0.736483
detractors        0.263517
Name: Sentiment, dtype: float64
After
0    3773
1    1350
Name: Sentiment, dtype: int64
0    0.736483
1    0.263517
Name: Sentiment, dtype: float64
```

- Visually Inspecting Features on train and test samples

```
In [36]: X_test.head()
```

```
Out[36]:
```

		Review	Rating	label
5372	great location good hotel 1night stay recommend hotel want just night accommodation ahead day flight, sheraton stars mentioned architecture n't good hotel, stay night try solution night stay choose directly connected airport, pretty expensive cola 0,2l juices 4 eur price minibar pizza 12 eur ok person problems finding signs sheraton, problems big frankfurt airport.actually easy hotel airport directly connected airport terminal 1.from departure hall upstairs 2 pedestrian bridges whichconnect terminal hotel long distance trains 1 hall b near lufthansa 1 hall b c near american airlines follow signs hotel fernbahnhof long distance trains,		4	not_detractors
16041	good airport hotel booked hotel early departures, nights room comfortable air-conditioning awful, quiet open windows, non-smoking floor smoking allowed restaurants common areas, arriving hotel shuttle available departing odd hour pay 30 taxi no taxi wants 5-minute trip, offer cab driver 10 did agree pay going rate hotel car 30, check-in problem odd hours staff speak english, hotel restaurants awful, chopsticks chinese taste good chinese food u.s. greasy, good wine list course problem smokers, good espresso desert bar outside chopsticks, best aspect hotel excellent travel desk manned efficient friendly lindsay day, ask eat town, english excellent, helped hotel merchant sold damaged item,		2	detractors

- Visually Inspecting Features on train and test samples

In [36]: x\_test.head()

Out[36]:

		Review	Rating	label
5372	great location good hotel 1night stay recommend hotel want just night accommodation ahead day flight, sheraton stars mentioned architecture n't good hotel, stay night try solution night stay choose directly connected airport, pretty expensive cola 0,2l juices 4 eur price minibar pizza 12 eur ok person problems finding signs sheraton, problems big frankfurt airport.actually easy hotel airport directly connected airport terminal 1.from departure hall upstairs 2 pedestrian bridges whichconnect terminal hotel long distance trains 1 hall b near lufthansa 1 hall b c near american airlines follow signs hotel fernbahnhof long distance trains,	4	not_detractors	
16041	good airport hotel booked hotel early departures, nights room comfortable air-conditioning awful, quiet open windows, non-smoking floor smoking allowed restaurants common areas, arriving hotel shuttle available departing odd hour pay 30 taxi no taxi wants 5-minute trip, offer cab driver 10 did agree pay going rate hotel car 30, check-in problem odd hours staff speak english, hotel restaurants awful, chopsticks chinese taste good chinese food u.s. greasy, good wine list course problem smokers, good espresso desert bar outside chopsticks, best aspect hotel excellent travel desk manned efficient friendly lindsay day, ask eat town, english excellent, helped hotel merchant sold damaged item,	2	detractors	
15225	mandarin oriental good stayed mandarin oriental way australia europe, visit hotel frankly difficult convince ourselevs stay singapore, stayed club rooms time service provided club guests second none, club rooms need little updating clean modern expect class hotel, location superb access marina square mall subsequently malls linked, prices risen significantly recent years watch special deals bonus free nights, highly recommended,	5	not_detractors	
18749	great little hotel close forbidden city great little hotel ideal location, minutes walking forbidden city, located traditional beijing hutong neighborhood recommend going early morning late night walk neighborhood sense locals live daily lives, flying shanghai beijing beijing airport got shut heavy rains flight canceled, thought hotel charged night no proper cancellation, arrived morning told not charge missed night knew control.sunny yu speaks good english help tours great wall places interested.my room bit small clean, overall great hotel grateful generously not charging missed night,	5	not_detractors	
8348	anther great stay husband stayed 9 nights july 2008. second stay not disappointed, location hotel fantastic rooms nicely appointment great sea view, short walk sitges centre quick taxi ride, make quiet beaches good restaurants 5-mins walking distance hotel.this hotel definitely gets busier weekends spanish tourists coming weekend break uncomfortable, hotel let slightly minor irritations, breakfast superb meals poor snack bar food terrible getting decent lunch hard, surprising considering brilliant restaurants area, needs improved, irritation lack uk tv channels, clearly does n't sitges watch tv day catch movies 9 days good, no uk channels cnn no movie ordering service, hotel clearly make money introducing service.all excellent choice sitges definitely return,	4	not_detractors	

```
In [37]: # Making a sample of 3 records to display the full text of each
train_sample = X_train.sample(3, random_state=22)
train_sample['target'] = [y_train[val] for val in train_sample.index]
train_sample.style.set_properties(**{'text-align': 'left'})
```

Out[37]:

		Review	Rating	label	target
		great location room able walk to/from train station no problems blocks away, location hotel excellent, kept map times easily able navigate way stopping ask directions time needed, people friendly, did morning guided walk florence recommend doing helped berings efficient way big sights david accademia gallery worth			
5452	seeing, took bus firenze nice view florence michelangelo piazza far	3	detractors	1	

```
In [38]: # Making a sample of 5 records to display the full text of each
exploring river, wonderful atmosphere, day night felt safe time
ended darker alleys accident no problems, food expensive nice
minus outdoor market smells leather coat my daughter loved fashion
test_sample = test.sample(5)
test_sample['target'] = test['target'].str.replace(' ', '_')
test_sample['target'].head()
for val in test_sample['target']:
    print(val)
test_sample.style.set_properties(**{'text-align': 'left'})
```

Out[38]:

		Review	Rating	label	target
	booking trip described paradisus excited time finally came, say overall hotel lived expectation 1. weather perfect entire time island.				
6365	didn't stay at night expected, temperature sun passed expectations, hot weekend he woyne island was great, bathroom luxurious quiet cool suite, course, mosquitoes not day, great service, food fast tea digital temperature control, beds comfy great sheets pillows, room leaving nice, did not major problems, hunts nothing major worked room morning, Weekend did not hear construction noises, 3rd floor no real view no street noise, bathroom huge ample supply luxurious towels, shower worked great good water pressure plenty hot water, rooms spotless, flat screen tv ipod dock nice touches, staff cordial helpful check check, wine hour 5pm 6pm real treat free coffee mornings, usually stay algonquin visits new york sold dates took chance muse, lucky not better choice.put place list not disappointed,	5	not_detractors	0	
330	no customer care booked travelocity, flight cancelled snow ice delayed day westin refused refund 1 2 nights, 200+ night hotel bit greedy, not mention view n' great air did n't work barely slept night, expect minimum air work room,still waiting hear travelocity not refund told left town refused got, happily ate cost rental can extra day.we wo		1	detractors	1

```
In [38]: # Making a sample of 1000 rows to display the full text of each
test_sample = X_test.sample(1000, random_state=42)
test_sample['label'].value_counts()
for val in test_sample.index:
    test_sample.style.set_properties(**{'text-align': 'left'})
    great hotel dominican republic booked trip paradisus punta cana
    march 2006 stayed hotel early september, websites reviewed prior
    booking trip described paradisus excited time finally came, say
    overall hotel lived expectation 1. weather perfect entire time island
    didn't rain not hot expected temperature double felt
    house stayed weekend sun 6-9 property surpassed expectations, hot
    weekend normal island weather daily rooms same day nice
    cool
    course, mosquitoes got day big red bites less rest trip did not
    comfortable, air conditioner worked perfectly quiet easy see digital
    temperature control, bathroom great new pillows fantastic
    nice did not major problems bugs ants nothing major worked room
    morning, weekend did not hear construction noises, 3rd floor no real
    6365 view no street noise, bathroom huge ample supply luxurious towels,
    shower worked great good water pressure plenty hot water, rooms
    spotless, flat screen tv ipod dock nice touches, staff cordial helpful
    check check, wine hour 5pm 6pm real treat free coffee mornings,
    usually stay algonquin visits new york sold dates took chance muse,
    lucky not better choice.put place list not disappointed,
    no customer care booked travelocity, flight cancelled snow ice
    delayed day westin refused refund 1 2 nights, 200+ night hotel bit
    greedy, not mention view n't great air did n't work barely slept night,
    expect minimum air work room.still waiting hear travelocity not refund
    told left town refused got, happily ate cost rental car extra day.we wo
    n't dealing time soon,
    great stay jesus berenguer lintag stayed ramada 4 days, service good
    16243 staff courteous helpful room amenities 4 star level, enjoyed stay wish
    say behalf wife daughter thank hotel staff, mabuhay,
```

	Review	Rating	label	target
6365	didn't rain not hot expected temperature double felt house stayed weekend sun 6-9 property surpassed expectations, hot weekend normal island weather daily rooms same day nice cool course, mosquitoes got day big red bites less rest trip did not comfortable, air conditioner worked perfectly quiet easy see digital temperature control, bathroom great new pillows fantastic nice did not major problems bugs ants nothing major worked room morning, weekend did not hear construction noises, 3rd floor no real view no street noise, bathroom huge ample supply luxurious towels, shower worked great good water pressure plenty hot water, rooms spotless, flat screen tv ipod dock nice touches, staff cordial helpful check check, wine hour 5pm 6pm real treat free coffee mornings, usually stay algonquin visits new york sold dates took chance muse, lucky not better choice.put place list not disappointed, no customer care booked travelocity, flight cancelled snow ice delayed day westin refused refund 1 2 nights, 200+ night hotel bit greedy, not mention view n't great air did n't work barely slept night, expect minimum air work room.still waiting hear travelocity not refund told left town refused got, happily ate cost rental car extra day.we wo n't dealing time soon,	5	not_detractors	0
16243	great stay jesus berenguer lintag stayed ramada 4 days, service good staff courteous helpful room amenities 4 star level, enjoyed stay wish say behalf wife daughter thank hotel staff, mabuhay,	4	not_detractors	0
330	great stay jesus berenguer lintag stayed ramada 4 days, service good staff courteous helpful room amenities 4 star level, enjoyed stay wish say behalf wife daughter thank hotel staff, mabuhay,	1	detractors	1

## 4: 2- Data Preprocessing & Exploratory Analysis

The dataset is being divided into two separate subsets: a training set, and a testing (or validation) set. The validation set will allow to assess the performance of the model.

Two parameters are assigned when dividing the dataset:

- random\_state=42
  - setting a random seed of 42 ensures that the data split is reproducible
- stratify=
  - stratified sampling ensures the class distribution is maintained in both sets to address potential class imbalance issues

In order to preprocess the reviews, the following transformations were performed:

- **Standardizing case**  
This step was verified as it is important to ensure text is uniform and consistent. Nevertheless, no extra step was taken as text was already saved as lower case.
- **Character Encoding**  
This was done to ensure consistent representation of text by transforming non-text characters into a normalized format, which will facilitate proper text processing, analysis, and model training.
- **Tokenizing**  
Tokens of words were created. This was done with the `RegexpTokenizer` package from `nltk.tokenize`. This step will facilitate the conversion of word into a suitable form for analysis and modeling.
- **Mutual Information Scores**
- **Stopwords**  
Before any transformation is done, a copy of the review column will be done so the original one can always be accessed.  
To focus on the data's theme, English stopwords were removed. Manual additions were made in this text's context, such as adding industry specific word (i.e. "hotel", "stay")

```
In [39]: # Duplicating the column review
X_train[original_review] = X_train['Review']
The WordNetLemmatizer package from nltk.stem.wordnet was used to reduce words to their base form, allowing a more accurate analysis
```

```
In [40]: # Verifying the new column was correctly created
The FreqDist package was used to review in a dictionary-like output, the words and their frequencies
```

```
Out[40]: • WordCloud
The words' frequencies were represented visually thanks to the WordCloud package
```

	Review	Rating	label	original_review
1282	gem, pleasantly surprised accommodations helpful attentive staff Bigraps were drawn to have a better understanding of the themes thanks to the bus street bit noisy did not away enjoyment nice boutique hotel, definitely stay time,	5	not_detractors	gem, pleasantly surprised accommodations helpful attentive staff Bigraps were drawn to have a better understanding of the themes thanks to the bus street bit noisy did not away enjoyment nice boutique hotel, definitely stay time,

- **Mutual Information Scores**  
and modeling.

Before any transformation is done, a copy of the review column will be done so the original one can always be accessed.  
To focus on the data's theme, English stopwords were removed. Manual additions were made in this text's context, such as adding industry specific word (i.e. "hotel", "stay")

In [39]:

```
# Duplicating the column review
# The WordNetLemmatizer package from nltk.stem.wordnet was used to reduce words to their
X_train[~X_train['Review']] = X_train['Review']
base form, allowing a more accurate analysis
```

In [40]:

```
# Verifying the new column was correctly created
# The Freqdist package was used to review in a dictionary-like output, the words and their
X_train[:1]
frequencies
```

Out[40]:

	Review	Rating	label	original_review
1	The words' frequencies were represented visually thanks to the WordCloud package			gem, pleasantly surprised
2	Bigrams	accommodations helpful attentive staff		accommodations helpful attentive staff
3	Bigrams	need, breakfast good cookies, room		need, breakfast good cookies, room
4	1282	clean comfortable, room fronting	Detractors	lovely clean comfortable, room fronting
5	bush street bit noisy did not away			bush street bit noisy did not away
6	enjoyment nice boutique hotel,			enjoyment nice boutique hotel,
7	definitely stay time,			definitely stay time,

## 4: 2-a) Standardizing Case

I will glance at the first sample of review to get an idea of whether I need to standardize case. This step is important to ensure text is standardized, to prevent models from treating words with different cases as different ones

In [41]:

```
# Isolating the first review into windows_sample
sample_review = train_sample.iloc[1]['Review']
sample_review
```

Out[41]:

"great hotel dominican republic booked trip paradisus punta cana march 2006 stayed hotel early september, websites reviewed prior booking trip described paradisus excited time finally came, say overall hotel lived expectation 1, weather perfect entire time island, did n't rain not not hot expected, temperature doab le felt comfortable normal island wear shorts skirts tank tops bathing suits co urse, mosquitoes got day big red bites legs rest trip did n't really itch, obviously recommend bring bug spray.2, rooms room nice, did n't major problems bugs ants nothing major, worked room nice large, bathroom probably needs facelift i t.3, hotel grounds amazing big, literally tram hotel, lucky room right pool are a great location, food pool beach right, just tram lobby area, grounds perfecti on, cute little huts benches sit walkways nice just relax dinner.4, food though t good, inclusives jamaica twice traveled parts world definitely rate food para disus, remember dominican republic not new york city quality n't really comparable, lunch buffet amazing, lobster want, n't believe, boyfriend probably ate 5 lobsters day lunch, wonderful, favorite place italian food japanese food, boyfr iend steak house brazilian food, el romanitico wonderful atmosphere went twice thought food better, 5. service did n't speak english thought, prepared spend p atient time explaining things special requests, thought service better grounds beach pool area restaurants good, 6. excursions did day trip island called sama ra, took 10 person plane area 40 minute trip bus tour horse riding trip mountai n hiked waterfall followed horses interested lunch freshly butchered chickens r ice boat trip island just coast buses plane ride, just 250 person, cheap got, c ompany called flying tours, 7. tipping tipped people thought deserved, simple,

In [42]:

```
# Checking that all reviews are in lowercase
def islowercase(df):
    df['islower'] = df['Review'].str.islower()
    print('The assumption that all reviews is + str(islower)')

islower = df['islower'].all()
print('The assumption that all reviews are in lowercase is ' + str(islower))
islower
```

It seemed clear that all the text was lowercase though. Let's filter on the reviews that were not, out of curiosity and extra verification.

In [43]:

```
Changing lowercase is not necessary, it looks like it was already done. I will verify whether this
is the case
special_characters_sample = X_train[~X_train['Review'].str.islower()][:2]
other_special_characters_sample = X_train[~X_train['Review'].str.islower()][10:12]
special_characters_sample
```

Out[43]:

	Review	Rating	label	original_review
1	time	riu	year	group travelling fifth trip
2	time	riu	year	group travelling fifth trip

```
In [42]: servers bartenders room keepers got money day, took 100 ones spent 60 people da  
ys helps plan, didn't think service better tipped did n't tip just did anyway  
def islowercase(df, column):  
    screen camera bug spray need, hotel charges fortune t  
imes, walk beach direction strisbeachside()all tourist shops items priced reasonably,  
dominican 'repubiic' known unique inca art, like said previous years stayed beach  
es negril jamaica, make comparison definitely choose negril dominican republic,  
hotel jamaica smaller Review  
The assumption that all reviews are in lowercase is False  
The assumption that all reviews are in lowercase is False  
paradisus place dominican republic culture want experience p  
paradisus excellent choice, enjoy goes quick,
```

It seemed clear that all the text was lowercase though. Let's filter on the reviews that were not, out of curiosity and extra verification.

```
In [43]: # Changing lowercase is not necessary, it looks like it was already done. Let's verify whether this  
# special_characters_sample = X_train[X_train['Review'].str.islower()[:2]  
# other_special_characters_sample = X_train[X_train['Review'].str.islower()][10:12]  
# special_characters_sample
```

Out[43]:

Review	Rating	label	original_review
time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i_Ã—_Ã©Ã—_ not staying ritz someplace like don_Ã—_Ã©_ expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomas mentioned asked pool view room based reviews i_Ã—_Ã©Ã—_Ã©_ read, score mostly 3 rooms 1 person enjoy	5	positive	time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i_Ã—_Ã©Ã—_ not staying ritz someplace like don_Ã—_Ã©_ expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomas mentioned asked pool view room based reviews i_Ã—_Ã©Ã—_Ã©_ read, score mostly 3 rooms 1 person enjoy

The uppercase text is actually characters that are incorrectly translated. I opened the file with the latin-1 encoding, let's try utf-8 as encoding the review if this resolves this.

```
In [44]: # Saving the index of the special_characters_sample to verify them in the future  
spec_char_indices = special_characters_sample.index
```

```
In [45]: # Loading dataset and saving it as raw_df  
utf8_df = pd.read_csv('data/tripadvisor_hotel_reviews.csv.zip', compression='zip')
```

```
In [46]: # Checking that all reviews are in Lowercase  
is_lowercase(utf8_df['Review'])
```

The assumption that all reviews are in lowercase is False

```
In [47]: # Inspecting 2 reviews that caused the reviews not to be Lowercase  
| utf8_df.loc[spec_char_indices]
```

**Out[47]:** Whether I open the file with encoding utf-8 or latin-1, I get the same results. I went in and checked the data source; it does contain those special characters from the source.

checked the data source. It does contain these special characters from the source.

time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i\_C\_é not staying ritz someplace like don\_C\_é expect resort provide service, fun regardless little quirks do.check in after 30 minutes bus ride resort arrived resort, 4: 2- b) Character Encoding Check out 10 people by resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day keys safe included package given session rooms I will then opt to keep the initial encoding mentioned asked pool view room based reviews C\_é e read group mainly 3 rooms next section ~~will encode characters which are 2011 and 1/2 so they dont affect the display~~ rooms, impression naiboa 1 type room turns rooms queen/double mixes time we\_C\_é try type room.we no problems water pressure hot water supply, power went couple times minute, 30 minutes afternoon long hassle.as mentioned air conditioning hotel tends weak non-existent rooms, rooms ceiling fan keeps room comfortable asking 4 times a/c gave up.i prepared solid beds weren\_C\_é nearly hard expected, understand rest group beds softer, said no problems sleeping, day sun drinks it\_C\_é easy sleep point concern room keys, mentioned room 2116 key opened 2117 2118 2119. knew rooms handy times big security risk, key engraved room 2170 i\_C\_é not sure work door not, point keys 2117 2118 2119 did not open door able open couple own.food think hotel offered widest selection breakfast items resort i\_C\_é e, granted selection problem getting morning started, think best thing french-toast style croissants yummy, breakfast terrace morning nice touch.lunch normally eaten beach pizza/beach bar/hamburgers/fries/pasta/burritos/roasted chicken good

Out[47]:

Whether I open the file with encoding `utf-8` or `latin-1`, I get the same results. I went in and checked the data source: it does contain these special characters from the source.

time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i\_\_é\_ | not staying ritz someplace like don\_\_é\_ expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, 4: 2- b) **Character Encoding**  
pleasenly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session roomas mentioned asked pool view room-based reviews |\_\_é\_ | read, group mainly 3 rooms next section will encode gharoters which are coming up to there don't affect modeling rooms, impression naiboa 1 type room turns rooms queen/double mixes time we\_\_é\_ try type room.we no problems water pressure hot water supply, power went couple times minute, 30 minutes afternoon long hassle.as mentioned air conditioning hotel tends weak non-existent rooms, rooms ceiling fan keeps room comfortable asking 4 times a/c gave up.i prepared solid beds weren\_\_é\_ nearly hard expected, understand rest group beds softer, said no problems sleeping, day sun drinks it\_\_é\_ easy sleep point concern room keys, mentioned room 2116 key opened 2117 2118 2119, knew rooms handy times big security risk, key engraved room 2170 i\_\_é\_ not sure work door not, point keys 2117 2118 2119 did not open door able open couple own.foodi think hotel offered widest selection breakfast items resort i\_\_é\_ e, granted selection problem getting morning started, think best thing french-toast style croissants yummy, breakfast terrace morning nice touch.lunch normally eaten beach pizzaerie beach bar burgers pizza pasta hot dogs, roasted chicken sand

6829

4

In [48]: # Reminding the previously created `special_characters_reviews`  
`special_characters_sample`

```
In [48]: # Reminding the previously created special_characters_reviews
special_characters_sample
```

Out[48]:

		Review	Rating	label	original_review
6829		time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i_Ã¢_Ã©Ã¢ not staying ritz someplace like don_Ã¶_Ã© expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomas mentioned asked pool view room based reviews i_Ã¢_Ã©Ã¢e read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms,	4	not_detractors	time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i_Ã¢_Ã©Ã¢ not staying ritz someplace like don_Ã¶_Ã© expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomas mentioned asked pool view room based reviews i_Ã¢_Ã©Ã¢e read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms,

Out[48]:

	Review	Rating	label	original_review
6829	<p>time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i_Ã¢_Ã©Ã¢ not staying ritz someplace like don_Ã¢_Ã© expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomsas mentioned asked pool view room based reviews i_Ã¢_Ã©Ã¢e read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms, impression naiboa 1 type room turns rooms queen/double mixes time we_Ã¢_Ã©Ã¶ try type room.we no problems water pressure hot water supply, power went couple times minute, 30 minutes afternoon long hassle.as mentioned air conditioning hotel tends weak non-existent rooms, rooms ceiling fan keeps room comfortable asking 4 times a/c gave up.i prepared solid beds weren_Ã¢_Ã© nearly hard expected, understand rest group beds softer, said no problems sleeping, day sun drinks it_Ã¢_Ã© easy sleep point concern room keys, mentioned room 2116 key opened 2117 2118 2119, knew rooms handy times big security risk, key engraved room 2170 i_Ã¢_Ã©Ã¢ not sure work door not, point keys 2117 2118 2119 did not open door able open couple own.foodi think hotel offered widest selection breakfast items resort i_Ã¢_Ã©Ã¢e, granted selection problem getting morning started, think best thing french-toast style croissants yummy, breakfast terrace morning nice touch.lunch normally eaten beach pizzeria beach bar burgers pizza fries pasta hot dogs, roasted chicken good swiss chalet, didn_Ã¢_Ã© try main buffet lunch don_Ã¢_Ã© know selection.supper different theme night buffet good, did italian la carte 3 times steak house caribbean restaurant, book a-la cartes need talk buffet manager night pick reservation slip breakfast day, snacks time couldn_Ã¢_Ã© eat, pool bar beach bar pizzeria food served 2, managed 24 hour food bamboo 5 minutes it_Ã¢_Ã© end supper buffet 2 seating_Ã¢_Ã© day arrive buffet manager book seating stay, assigned</p>	4	not_detractors	<p>time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i_Ã¢_Ã©Ã¢ not staying ritz someplace like don_Ã¢_Ã© expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomsas mentioned asked pool view room based reviews i_Ã¢_Ã©Ã¢e read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms, impression naiboa 1 type room turns rooms queen/double mixes time we_Ã¢_Ã©Ã¶ try type room.we no problems water pressure hot water supply, power went couple times minute, 30 minutes afternoon long hassle.as mentioned air conditioning hotel tends weak non-existent rooms, rooms ceiling fan keeps room comfortable asking 4 times a/c gave up.i prepared solid beds weren_Ã¢_Ã© nearly hard expected, understand rest group beds softer, said no problems sleeping, day sun drinks it_Ã¢_Ã© easy sleep point concern room keys, mentioned room 2116 key opened 2117 2118 2119, knew rooms handy times big security risk, key engraved room 2170 i_Ã¢_Ã©Ã¢ not sure work door not, point keys 2117 2118 2119 did not open door able open couple own.foodi think hotel offered widest selection breakfast items resort i_Ã¢_Ã©Ã¢e, granted selection problem getting morning started, think best thing french-toast style croissants yummy, breakfast terrace morning nice touch.lunch normally eaten beach pizzeria beach bar burgers pizza fries pasta hot dogs, roasted chicken good swiss chalet, didn_Ã¢_Ã© try main buffet lunch don_Ã¢_Ã© know selection.supper different theme night buffet good, did italian la carte 3 times steak house caribbean restaurant, book a-la cartes need talk buffet manager night pick reservation slip breakfast day, snacks time couldn_Ã¢_Ã© eat, pool bar beach bar pizzeria food served 2, managed 24 hour food bamboo 5 minutes it_Ã¢_Ã© end supper buffet 2 seating_Ã¢_Ã© day arrive buffet manager book seating stay, assigned</p>
3320	<p>monday night saw first baseball game thursday evenings admitted saw people turned street bars set street vendors, pool bar beach bar lobby bar night guests involved actinether animation team chamber plus did come beach spot more guests remind events 2009 participants came pleased stand basically white, bent over, expect need watch stretching sessions water volleyball walk smokers adjust sitting water polo chose, ashtrays help habit.we spent 1 caribbean beach chairs around poolside finding available nights see second apprehensive about reviews arrived beach checkers give grueling loggers gave, stayed dry, heard room noise nail we put action movie screen new beach neighbor a hotel see night wasnt water full load control generally spend day beach want reading restaurants old fashioned brittle baity using their entertainment area above audience participation people really fun, expensive, unique, interesting, great work awesome, some review better,</p>	4	not_detractors	<p>monday night saw first baseball game thursday evenings admitted saw people turned street bars set street vendors, pool bar beach bar lobby bar night guests involved actinether animation team chamber plus did come beach spot more guests remind events 2009 participants came pleased stand basically white, bent over, expect need watch stretching sessions water volleyball walk smokers adjust sitting water polo chose, ashtrays help habit.we spent 1 caribbean beach chairs around poolside finding available nights see second apprehensive about reviews arrived beach checkers give grueling loggers gave, stayed dry, heard room noise nail we put action movie screen new beach neighbor a hotel see night wasnt water full load control generally spend day beach want reading restaurants old fashioned brittle baity using their entertainment area above audience participation people really fun, expensive, unique, interesting, great work awesome, some review better,</p>

	Review	Rating	label	seating_AU_AC day arrive buffet manager book seating stay, assigned	original_review
3320	specific little stay caribbean beach thursday night caravani turned turned away bus there since was worse; pool bar beach parades street gets guests 7/17/06, and animators tai no bamboo plus didn't come beach spending good time events bus party, punta cana pleased stand basically white stretching sessions water volleyball walk smokers adjust putting puts sand chose, ashtrays help habit we spent 1 caribean beach 7/17/06, and animators finding new hotel nights june 2005; apprehensive about received and beach checkers give great loggers gave couple try hard to do naihan pool action across new beach, not hotel A hotel swimming waste water not bad continue more ready spend play beach entertainment restauraants band live on the beach using their entertainment was about audience participation provide really fun, explosive, electric, breaking beat work awesome, especially better	4	not_detractors	specific little stay caribbean beach thursday night caravani turned turned away bus there since was worse; pool bar beach parades street gets guests 7/17/06, and animators tai no bamboo plus didn't come beach spending good time events bus party, punta cana pleased stand basically white stretching sessions water volleyball walk smokers adjust putting puts sand chose, ashtrays help habit we spent 1 caribean beach 7/17/06, and animators finding new hotel nights june 2005; apprehensive about received and beach checkers give great loggers gave couple try hard to do naihan pool action across new beach, not hotel A hotel swimming waste water not bad continue more ready spend play beach entertainment restauraants band live on the beach using their entertainment was about audience participation provide really fun, explosive, electric, breaking beat work awesome, especially better	seating_AU_AC day arrive buffet manager book seating stay, assigned

I will try to encode these characters by normalizing the Unicode string using NFKD (Normalization Form Compatibility Decomposition). After the string is normalized, they are encoded to ASCII and any characters that cannot be represented will be ignored, and then it is encoded back to latin-1.

In [49]: package

```
nction to normalize the characters.  
e['Review'] = special_characters_sample['Review'].apply(lambda x: unicodedata.nor
```

In [50]: # Reviewing the modified review  
special\_characters\_sample

Out[50]:

	Review	Rating	label	original_review
	time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i_A_AA not staying ritz someplace like don_A_A expect resort provide service, fun regardless little quirks do.check in after 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.rooms mentioned asked pool view room based reviews i_A_AAE read,			time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i_A_AAC not staying ritz someplace like don_A_AAC expect resort provide service, fun regardless little quirks do.check in after 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.rooms mentioned asked pool view room based reviews i_A_AAE read,

In [51]: special\_characters\_sample['Review'] = special\_characters\_sample['Review'].str.replace('Ã', 'A')  
special\_characters\_sample['Review'] = special\_characters\_sample['Review'].str.replace('Ã©', 'E')  
This helped! Now, our special characters are replaced by \_A\_. Will remove them from the  
reviews on the next cell.  
# Verifying it removed it  
special\_characters\_sample

```
In [51]: special_characters_sample['Review'] = special_characters_sample['Review'].str.replace('[^A-Za-z]', '_AA_')  
# Verifying it removed it  
special_characters_sample
```

Out[51]:

		Review	Rating	label	original_review
6829		time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i not staying ritz someplace like don expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomsas mentioned asked pool view room based reviews ie read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms, impression naiboa 1 type room turns rooms queen/double	4	not_detractors	time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i ___. expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomsas mentioned asked pool view room based reviews ie read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms,

Out[51]:

		Review	Rating	label	original_review
6829		<p>time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i not staying ritz someplace like don expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.rooms mentioned asked pool view room based reviews ie read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms, impression naiboa 1 type room turns rooms queen/double mixes time we try type room.we no problems water pressure hot water supply, power went couple times minute, 30 minutes afternoon long hassle.as mentioned air conditioning hotel tends weak non-existent rooms, rooms ceiling fan keeps room comfortable asking 4 times a/c gave up.i prepared solid beds weren nearly hard expected, understand rest group beds softer, said no problems sleeping, day sun drinks it easy sleep point concern room keys, mentioned room 2116 key opened 2117 2118 2119, knew rooms handy times big security risk, key engraved room 2170 i not sure work door not, point keys 2117 2118 2119 did not open door able open couple own.foodi think hotel offered widest selection breakfast items resort ie, granted selection problem getting morning started, think best thing french-toast style croissants yummy, breakfast terrace morning nice touch.lunch normally eaten beach pizzeria beach bar burgers pizza fries pasta hot dogs, roasted chicken good swiss chalet, didn't try main buffet lunch don know selection.supper different theme night buffet good, did italian la carte 3 times steak house caribbean restaurant, book a-la cartes need talk buffet manager night pick reservation slip breakfast day, snacks time couldn't eat, pool bar beach bar pizzeria food served 2, managed 24 hour food bamboo 5 minutes it end supper buffet 2 seating, day arrive buffet manager book seating stay, assigned specific table stay, men wear pants shirt sleeves admitted saw people turned away.bars there bars visit wish, pool bar beach bar lobby bar night pizzeria bar bamboo bar swimwear</p>	4	not_detractors	<p>time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i not staying ritz someplace like don expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.rooms mentioned asked pool view room based reviews ie read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms, impression naiboa 1 type room turns rooms queen/double mixes time we try type room.we no problems water pressure hot water supply, power went couple times minute, 30 minutes afternoon long hassle.as mentioned air conditioning hotel tends weak non-existent rooms, rooms ceiling fan keeps room comfortable asking 4 times a/c gave up.i prepared solid beds weren nearly hard expected, understand rest group beds softer, said no problems sleeping, day sun drinks it easy sleep point concern room keys, mentioned room 2116 key opened 2117 2118 2119, knew rooms handy times big security risk, key engraved room 2170 i not sure work door not, point keys 2117 2118 2119 did not open door able open couple own.foodi think hotel offered widest selection breakfast items resort ie, granted selection problem getting morning started, think best thing french-toast style croissants yummy, breakfast terrace morning nice touch.lunch normally eaten beach pizzeria beach bar burgers pizza fries pasta hot dogs, roasted chicken good swiss chalet, didn't try main buffet lunch don know selection.supper different theme night buffet good, did italian la carte 3 times steak house caribbean restaurant, book a-la cartes need talk buffet manager night pick reservation slip breakfast day, snacks time couldn't eat, pool bar beach bar pizzeria food served 2, managed 24 hour food bamboo 5 minutes it end supper buffet 2 seating, day arrive buffet manager book seating stay, assigned specific table stay, men wear pants shirt sleeves admitted saw people turned away.bars there bars visit wish, pool bar beach bar lobby bar night pizzeria bar bamboo bar swimwear</p>
		<p>arrived at noon plus 15 min delay come beach spot hotel guest friendly car pleasant and has great service beach bar open early need catch walk voicemail adult smoking parts said chose, ashtrays help habit.we spent 1 day beach, chairs want sit problem finding available ones use second week quite crowded paid beach workers reserve group 8 loungers easy couple trying 8 heard, poolwe didn't pool got ones place stayed beach bar into night was 2006 about end of day never play dress was an old set of regal maya asked upgrade gave, stay getting value quite new to them great team involve 2 new, beach use performance sove right time where we much less jackson day work average some relaxation being movements like vacation feel deal quite early every night review vr maintained cabanas bassed at vr vendor highlight quiet paradise strait gear quest reviewed activities the</p>	4	not_detractors	<p>arrived at noon plus 15 min delay come beach spot hotel guest friendly car pleasant and has great service beach bar open early need catch walk voicemail adult smoking parts said chose, ashtrays help habit.we spent 1 day beach, chairs want sit problem finding available ones use second week quite crowded paid beach workers reserve group 8 loungers easy couple trying 8 heard, poolwe didn't pool got ones place stayed beach bar into night was 2006 about end of day never play dress was an old set of regal maya asked upgrade gave, stay getting value quite new to them great team involve 2 new, beach use performance sove right time where we much less jackson day work average some relaxation being movements like vacation feel deal quite early every night review vr maintained cabanas bassed at vr vendor highlight quiet paradise strait gear quest reviewed activities the</p>
3320		<p>tower great ocean view beach, beach club resort hotel schoolhouse wasnt water apt had choice more than spend play game room including afternoon tea busan getting chairs easy using their service room by 2nd floor audience participation people really fun especially michel jackson foot work awesome some selection better</p>	4	not_detractors	<p>tower great ocean view beach, beach club resort hotel schoolhouse wasnt water apt had choice more than spend play game room including afternoon tea busan getting chairs easy using their service room by 2nd floor audience participation people really fun especially michel jackson foot work awesome some selection better</p>

	Review	Rating	label	seating_AU_AO_, day arrive butter manager book seating stay assigned specific table stay night beach shirts sleeves admitted raw people turned away bags storage space Mosh was pool bar beach bar lobby bar night high beach paradise street gets guest friendly admiring the animation tape hampon plus digging AU_AO bar team activities going day come beach were more guests behind them who party CAB pleased And basically whole beach clear XMAS red wall wall stretching sessions rate volleyball smokers just putting this same those ashtrays help habit we spent 1 day cabbed in towels want sit problem finding available car seats no longer work 2005 apartment night reviews star review group bookings easy asked bying 8 heart puzzle renovation look new tower great recent happy beach beach deck chairs hotel swimming waste water A lot had come more ready spend day beach room more date afternoon night than sitting chairs outside using the end of pier when about audience participation provide really fun expansive比如backstagefeat work awesome experience better
3320	an amazing place also good very good drink the beach nice come beach spot the terrace is the even pleasure and has a nice menu beach clean especie need watch walk volleyball and basketball pool chose, ashtrays help habit we spent 1 day beach, chairs want sit problem finding available ones use second week quite crowded paid beach workers reserve group 8 loungers easy couple trying 8 heard, poolwe decoration park and places stayed beach and the night was good work aware some entertainment movement like the main entrance easy entry easy reverse nailhead castles expensive but very highlight and provide start great guests provided recently after	4	not_detractors	

Great, this resolved it. I will save the steps into a function and will apply it to the full train data

In [52]: # Creating the function character\_encoding to handle them

```
def character_encoding(df_column):
    # Importing the relevant package
    import unicodedata

    # Normalizing characters
    df_column = df_column.apply(lambda x: unicodedata.normalize('NFKD', x).encode('utf-8'))

    # Generalizing to list of characters to remove
    characters_to_remove = ['_A_AA', '_A_A', '_AA', '_A', '_Aow_', 'AAA', '_AAo', '_AOA']

    # Removing the characters from text
    for char in characters_to_remove:
        df_column = df_column.str.replace(char, '')

    # Returning the df_column
    return df_column
```

In [53]: # Testing the function  
character\_encoding(other\_special\_characters\_sample['Review'])

Out[53]: 4875

wow oh wow, partner booked surprise, stayed fabulous room poster bed wonderfull y luxurious feel, think tudor suite, wished winter lit sat snug area, terry but ler great guide hotel local area, entertaining informative, went bar pre-dinner drinks daniel barman welcoming stephan offered champagne champagne day day mile stone, just fantastic lovely nibbles happy stayed bar night, dinner cheneston w onderful experience lovely setting, felt important special not just partner tre ated, zana lovely treated like royalty did staff, meal truly treat, started end ive gruyere pancetta terrine delicious main course, instead went pot roast dove r sole world, unfortunately n't manage dessert did steal mouthful simon souffl oh goodness, thomas sure excellent wines course convinced cognac coffee lounge, used, breakfast calm quality food service near perfect, can thank staff making short break memorable special ca n't loved milestone ca n't wait return just ar range girls night dinner, sure stephan thomas daniel make great night,  
19458 best vacation spotever best vacation, went large group, imagine travel ing large group cumbersome relief, need speak little spanish listen carefully e nglish it totally doable repetition, went way accommodate, eager make stay exce lente, tips, n't expecting seasons know surroundings, 3rd world country, beds k inda firm cater europeans americans little flexibile days change linens everuda

```
In [53]: # Testing the function
character_encoding(other_special_characters_sample['Review'])
```

Out[53]: 4875

wow oh wow, partner booked surprise, stayed fabulous room poster bed wonderfull y luxurious feel, think tudor suite, wished winter lit sat snug area, terry but ler great guide hotel local area, entertaining informative, went bar pre-dinner drinks daniel barman welcoming stephan offered champagne champagne day day mile stone, just fantastic lovely nibbles happy stayed bar night, dinner cheneston w onderful experience lovely setting, felt important special not just partner tre ated, zana lovely treated like royalty did staff, meal truly treat, started end ive gruyere pancetta terrine delicious main course, instead went pot roast dove r sole world, unfortunately n't manage dessert did steal mouthful simon souffl oh goodness, thomas sure excellent wines course convinced cognac coffee lounge, used, breakfast calm quality food service near perfect, can thank staff making short break memorable special ca n't loved milestone ca n't wait return just ar range girls night dinner, sure stephan thomas daniel make great night,  
19458 best vacation spotever best vacation, went large group, imagine travel ing large group cumbersome relief, need speak little spanish listen carefully e nglish it totally doable repetition, went way accommodate, eager make stay exce lente, tips, n't expecting seasons know surroundings, 3rd world country, beds k inda firm cater europeans americans little flexible days, change linens everyda y refresh liquor room days, food water restaurants safe, n't believe things wri tten, just got week ago not sick, food fresh lacks salt prefer, totally purifie d gallon water stocked fridge, drinking doing drink water, rarely saw individua l bottles water, ice make drinks water serve purified n't believe hype, beach t otally open safe night, careful waves undertow night, fierce, love juan carlos froggy, hotel, total sweethearts, went club pacha, club rocked, site shuttle no cabs leave safety compound, oh yeah did mention pools heated poolside bar, best thing swimming cold ocean, shopping site not expensive considering it thing vac ation pay, not sure place conducive kiddos available booze alcohol all-inclusiv e thing winner, not mention stay riu wristband gives access area hotels, make s ure breakfast sign dinner italian restaurant brazilian steakhouse days, fabulou s, sign 9-12pm dinner night, try sea bass italian place brazilian restaurant ba con wrapped bananas, sounds gross trust die, brazilian spot there lot meat litt le veggies prepared, tipped lot staff great, couple dollars long way, aware sta ff flirts, smooth totally caught, watch, riu punta cana, best vaca say hi georg ina jefe eugenio definitely possibly yearly,

Name: Review, dtype: object

I could now apply it to our full X\_train['Review'] column directly, as I kept an unmodified copy as the original. I later defined a function combining all preprocessing steps, so the next cells will be commented out.

```
In [54]: # Reviewing the other sample
# other_special_characters_sample
```

```
In [55]: # Calling the function on full X_train
# X_train['Review'] = character_encoding(X_train['Review'])
```

```
In [56]: # X_train[~X_train['Review'].str.islower()]
```

```
In [57]: # Verifying that words containing 'a' lowercase were not removed
# assert len(X_train[X_train['Review'].str.contains('stay')]) > 0
```

```
In [58]: # Reviewing one of our sample of 2 reviews
# Verifying that our sample reviews are correctly handled now
special_characters_sample
# X_train.loc[spec_char_indices]
```

```
In [59]: # Verifying that I correctly have lowercase everywhere
# is_lowercase(X_train['Review'])
```

This now resolved it and gave me the result I hoped for. I can move on to the next step.

## 4: 2- c) Tokenizing

Tokenizing text data is one of the fundamental data cleaning steps to further convert words into a suitable form for analysis and modeling.

```
In [68]: # Reviewing one of our sample of 2 reviews  
In [69]: ## Verifying that our sample reviews are correctly handled now  
## special characters sample  
## X_train.loc[spec_char_indices]
```

```
In [59]: # Verifying that I correctly have Lowercase everywhere  
# is_lowercase(X_train['Review'])
```

This now resolved it and gave me the result I hoped for. I can move on to the next step.

## 4: 2- c) Tokenizing

Tokenizing text data is one of the fundamental data cleaning steps to further convert words into a suitable form for analysis and modeling.

Out[60]:

Out[60]:

		Review	Rating	label	original_review
6829		time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i not staying ritz someplace like don expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomsas mentioned asked pool view room based reviews ie read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms, impression naiboa 1 type room turns rooms queen/double mixes time we try type room.we no problems water pressure hot water supply, power went couple times minute, 30 minutes afternoon long hassle.as mentioned air conditioning hotel tends weak non-existent rooms, rooms ceiling fan keeps room comfortable asking 4 times a/c gave up.i prepared solid beds weren nearly hard expected, understand rest group beds softer, said no problems sleeping, day sun drinks it easy sleep point concern room keys, mentioned room 2116 key opened 2117 2118 2119. knew rooms handy times big security risk, key engraved room 2170 i not sure work door not, point keys 2117 2118 2119 did not open door able open couple own.foodi think hotel offered widest selection breakfast items resort ie, granted selection problem getting morning started, think best thing french-toast style croissants yummy, breakfast terrace morning nice touch.lunch normally eaten beach pizzeria beach bar burgers pizza fries pasta hot dogs, roasted chicken good swiss chalet, didn't try main buffet lunch don know selection.supper different theme night buffet good, did italian la carte 3 times steak house caribbean restaurant, book a-la cartes need talk buffet manager night pick reservation slip breakfast day, snacks time couldn't eat, pool bar beach bar pizzeria food served 2, managed 24 hour food bamboo 5 minutes it end supper buffet 2 seating, day arrive buffet manager book seating stay, assigned specific table stay, men wear pants shirt sleeves admitted saw people turned away.barsthere bars visit wish, pool bar beach bar lobby bar night pizzeria bar bamboo bar swim bar an animation show service go away come beach spot hotel guest time anima events poor party don't have all the beach clean experie need watch water volleyball and things water polo chose, ashtrays help habit.we spent 1 day beach, chairs want sit problem finding available ones use second week quite crowded paid beach workers reserve group 8 loungers easy couple trying 8 heard, poolwe callention pot and places shaped beach chair and pool was added after morning bad review day arrived want sit down again forward siting upgrade gave, stay better than room new when great time we can view, beach area perdition several day time where it much not go away work aware some selection better movement like other than paid qaré easy today chit reviewers narrated castries has been the verdays highlights such provide strength guest, perched polyester,	4	not_detractors	time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i _Ã_Ã© not staying ritz someplace like don _Ã_Ã© expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.roomsas mentioned asked pool view room based reviews i _Ã_Ã© read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms, impression naiboa 1 type room turns rooms queen/double mixes time we _Ã_Ã© try type room.we no problems water pressure hot water supply, power went couple times minute, 30 minutes afternoon long hassle.as mentioned air conditioning hotel tends weak non-existent rooms, rooms ceiling fan keeps room comfortable asking 4 times a/c gave up.i prepared solid beds weren _Ã_Ã© nearly hard expected, understand rest group beds softer, said no problems sleeping, day sun drinks it _Ã_Ã© easy sleep point concern room keys, mentioned room 2116 key opened 2117 2118 2119. knew rooms handy times big security risk, key engraved room 2170 i _Ã_Ã© not sure work door not, point keys 2117 2118 2119 did not open door able open couple own.foodi think hotel offered widest selection breakfast items resort i _Ã_Ã©, granted selection problem getting morning started, think best thing french-toast style croissants yummy, breakfast terrace morning nice touch.lunch normally eaten beach pizzeria beach bar burgers pizza fries pasta hot dogs, roasted chicken good swiss chalet, didn't _Ã_Ã© try main buffet lunch don _Ã_Ã© know selection.supper different theme night buffet good, did italian la carte 3 times steak house caribbean restaurant, book a-la cartes need talk buffet manager night pick reservation slip breakfast day, snacks time couldn't _Ã_Ã© eat, pool bar beach bar pizzeria food served 2, managed 24 hour food bamboo 5 minutes it _Ã_Ã© end supper buffet 2 seating _Ã_Ã© day arrive buffet manager book seating stay, assigned shirt sleeves admitted saw people turned away.barsthere bars visit wish, pool bar beach bar lobby bar night pizzeria bar bamboo bar swim bar an animation show service go away come beach spot hotel guest time anima events poor party don't have all the beach clean experie need watch water volleyball and things water polo chose, ashtrays help habit.we spent 1 day beach, chairs want sit problem finding available ones use second week quite crowded paid beach workers reserve group 8 loungers easy couple trying 8 heard, poolwe callention pot and places shaped beach chair and pool was added after morning bad review day arrived want sit down again forward siting upgrade gave, stay better than room new when great time we can view, beach area perdition several day time where it much not go away work aware some selection better movement like other than paid qaré easy today chit reviewers narrated castries has been the verdays highlights such provide strength guest, perched polyester,
3320				4	not_detractors

	pool bar beach bar lobby bar night pizzeria bar bamboo bar sun room	Review	Rating	label
3320	animat team activities going away comped leading to other dests elmina events poor park and basic beach field stretching sessions work smokers adjust sitting butts water volleyball aqua fitness water pond chose, ashtrays help habit.we spent 1 day beach, chairs want sit problem finding available ones use second week quite crowded paid beach workers reserve group 8 loungers easy couple trying 8 heard, poolwe caribean not pot ones places have beach chair problem was 2100 ashtray morning has new day beach water and lots of time we had upgrade gave, stay get the value route new toher great beautiful we can review, beach fence postination sove each A full view water sun nishon lack go away work average, some selection better movement and its own garden gone early today eat, cover every narrabeach cafrees bars restaurants venuelty highlight not made strait get guess as jeweled lobby hotel,	Review	4	not_detractors
	seating _A@_A@_d, arrive butte manager book seating stay review spec table stay in wear pants movement like Watch rear beach thursday night caravans can be seen turned away, bars see street vissib pro high Patrick barbeque streetside gues to drive, achim the swimwear tang harmons usd day come beach went open air bars the beach via gate spend time, guest seat, sunbath party @2000 please A and basicly with hold beach clean session, need water yellow smokers just passing fitness water pool ashtrays help habit.we spent 1 day change in hair, A at sit poleole finaly available ames inter friend week 2001 apprentice had as achievement reser group 2001, our survey ask to you he give, review is not about go there tower great beach, Japen beach, beach dolphin bar hotel pool right waste water A for bad choice go away spend day beach with holding afternoon hand can getting chairs easily using the service to go to review audience as explosive provide really fun, playfully, Michael working feel work awesome, exponed to even better,	Review	4	neutral_review

I will use `RegexpTokenizer` from NLTK to create tokens of two or more consecutive word characters, which include letters, numbers and underscores.

- Tokenizing Pattern

The token pattern was defined as follows:

- `r` indicates the string is in Python language
  - `(?u)` allows to match Unicode characters
  - `\b` determines where the word starts or ends
  - `\w` matches any word character
  - `\w+` matches 1 or more word characters
  - `\b` finally the word boundary ensures the matched pattern ends at a word boundary

```
In [61]: # Importing RegexpTokenizer
from nltk.tokenize import RegexpTokenizer

token_pattern = r"(?u)\b\w+\b"

# Instantiating the tokenizer
tokenizer = RegexpTokenizer(token_pattern)

# Tokenizing the reviews
special_characters_sample['Review'] = special_characters_sample['Review'].apply(
    lambda x: tokenizer.tokenize(x))
```

```
Out[61]: 6829 [time, riu, year, group, travelling, fifth, trip, dominican, republic, visit, punta, cana, travel, understanding, not, staying, ritz, someplace, like, don, expect, resort, provide, service, fun, regardless, little, quirks, do, check, inafter, 30, minute, bus, ride, resort, arrived, resort, check, good, 10, p
```

```
In [62]: # Importing the needed full tokenization library
# Xtrain['Review'].apply(lambda x: tokenizer.tokenize(x))
```

```
In [63]: days, keepingafe, included, package, given, session, rooms, mentioned, asked, p  
pol, tvæw, hoom) based, reviews, ie, read, group, ...]
```

3320 [caribe, hilton, piAA\_a, coladas, stayed, caribe, hilton, nights, june, 2005, a  
4: 2-d Stopwords  
pprehensive, bad, reviews, arrived, staff, checked, right, away, asked, upgrade,  
e, gave, stayed, renovated, room, new, tower, great, beautiful, ocean, view, be-  
ach, clean, order, hotel, serve, right, slide, water, lot, fun, concierge, real  
Then I will be removing stopwords and can focus on the text data element. It looked like the  
text was somewhat preprocessed so will review whether there is a need to further extract  
stopwords.  
Name: Review, dtype: object

This step will not be applied to the full `X_train` data as I will use it separately to test whether it positively impacts our model in the future when I reach this step.

I later defined a function combining all preprocessing steps, so the next cells will be commented out.

In [64]: # Importing relevant packages  
import altair

```
In [62]: # Applying the tokenizer on the full train data
# Xtrain['Review'] = Xtrain['Review'].apply(lambda x: tokenizer.tokenize(x))

In [63]: days, keepingafe, included, package, given, session, roomas, mentioned, asked, p
          pol, tvæw, head) based, reviews, ie, read, group, ...
3320
```

[caribe, hilton, piAA\_a, coladas, stayed, caribe, hilton, nights, june, 2005, apprehensive, bad, reviews, arrived, staff, checked, right, away, asked, upgrade, gave, stayed, renovated, room, new, tower, great, beautiful, ocean, view, beach, clean, order, hotel, serve, right, slide, water, lot, fun, concierge, real helpful, removing stopwords, cautious on the text data since it looks like the text was somewhat preprocessed, so will review whether there is a need to further extract Stopwords  
Name: Review, dtype: object

I could now apply it to our full X\_train data as I will use it separately to test whether it positively impacts our model in the future when I reach this step.

I later defined a function combining all preprocessing steps, so the next cells will be commented.

```
In [64]: # Importing relevant package
```

```
import nltk
```

```
nltk.download('stopwords', quiet=True)
```

```
from nltk.corpus import stopwords
```

```
# Creating list to store stopwords  
stopwords_list = stopwords.words('english')
```

```
Out[64]: ['i', 'me', 'my', 'myself', 'we']
```

I will continue using our previously created special\_characters\_sample

```
In [65]: special_characters_sample
```

Out[65]:

Review	Rating	label	original_review
[time, riu, year, group, travelling, fifth, trip, dominican,	5	positive	time riu year group travelling fifth trip dominican republic visit punta cana, travel understanding i_Ã¢_Ã© not staying ritz someplace like don_Ã¢_Ã© expect resort provide service, fun regardless little quirks do.check inafter 30 minute bus ride resort arrived resort, check good, 10 people bus resort 7 group, emailed week asking pool view room upgrade pleasantly surprised gotten, walked lobby naturally headed reception desk called small table given registration form room keys brought meeting room cocktail drink given general info hotel hotel not tour operator day, keys safe included package given session.rooms mentioned asked pool view room based reviews i_Ã¢_Ã© read, group mainly 3 rooms 1 person away group room issue, rooms 2116 2117 2118.there 1 double bed 1 single bed rooms, impression naiboa 1 type room turns rooms queen/double mixes time we_Ã¢_Ã© try type room.we no problems water pressure hot water supply. power went couple times

```
In [66]: # Defining function that takes in a List of strings and returns only those that a
```

```
def remove_stopwords(token_list, stopwords_list):
```

```
In [68]: # Reviewing token example review - stopwords removed = [token for token in token_list if token not in stopwords]
tokens_example_without_stopwords[:5]
```

Out[68]: ['time', 'riu', 'year', 'group', 'travelling']

```
In [67]: # Testing it on the tokens example
I will review during the exploratory analysis whether some words need to be manually added to the
tokens_example = special_characters_sample.iloc[6][Review]
list
print("Length with stopwords: ", len(tokens_example))
```

```
tokens_example_without_stopwords = remove_stopwords(tokens_example, stopwords_list)
print("Length without stopwords: ", len(tokens_example_without_stopwords))
```

Length with stopwords: 526  
The WordNet stopword package from nltk.stem.wordnet was used to reduce words to their base

It first required to be downloaded for Jupyter Notebook. Once the initial download is done, this step was commented out - the main reasoning is loss of risk of being impacted.

This step will not be applied to the full `X_train` data, as I will use it separately to test whether it positively impacts our model in the future when I reach this step.

```
In [68]: def remove_stopwords(token_list, stopwords_list):
    # Reviewing a token example
    tokens_removed = []
    for token in token_list:
        if token not in stopwords_list:
            tokens_removed.append(token)
    return tokens_removed

Out[68]: ['time', 'riu', 'year', 'group', 'travelling']

In [67]: # Testing it on the tokens example
    tokens_example = special_characters_sample.iloc[0][Review]
    print("Length with stopwords: ", len(tokens_example))

    tokens_example_without_stopwords = remove_stopwords(tokens_example, stopwords_list)
    print("Length without stopwords: ", len(tokens_example_without_stopwords))
```

Length with stopwords: 526  
The WordNetLemmatizer package from nltk.stem.wordnet was used to reduce words to their base form, allowing a more accurate analysis.  
It's required to be downloaded for Jupyter Notebook. Once the initial download is done, this step was commented out.

This step will not be applied to the full X\_train data, as I will use it separately to test whether it positively impacts our model in the future when I reach this step.

```
In [69]: # Importing relevant package
from nltk.stem.wordnet import WordNetLemmatizer
# nltk.download('wordnet')
# nltk.download('omw-1.4')

# Instantiating the Lemmatizer
lemmatizer = WordNetLemmatizer()

In [70]: # Instantiating the Lemmatizer
def lemmatize_words(token_list):
    lemmatized_tokens = [lemmatizer.lemmatize(token, pos='v') for token in token_list]
    return lemmatized_tokens

In [71]: # Print the tokens_example without stopwords before
print('Before')
print(tokens_example_without_stopwords[:5])

Before
['time', 'riu', 'year', 'group', 'travelling']

In [72]: # Let's review the impact on our token example
print('After')
lemmatize_words(tokens_example_without_stopwords)[:5]

After
['time', 'riu', 'year', 'group', 'travel']
```

The words were lemmatized as intended and reduced to their base form. This can be visible on words: travelling changed into travel.

#### 4: 2-f) Summarizing Review Preprocessing

```
In [74]: # Defining the function
def preprocess_review(df_column):
    # Previous steps will be gathered into one main function so the above steps can be applied to the entire dataset
    df_column = df_column.str.lower()

In [73]: # Ensuring character encoding
# Calling the previously created function
# These were imported individually before but are reminded here if they needed to be
df_column = character_encoding(df_column)
# from nltk.tokenize import RegexpTokenizer
# nltk.download('stopwords', quiet=True)
# from nltk.corpus import stopwords
# from nltk.stem.wordnet import WordNetLemmatizer
token_pattern = r'(?u)\b\w+\b'
# Instantiating the tokenizer
tokenizer = RegexpTokenizer(token_pattern)
# Tokenizing
# df_column = tokenizer.tokenize(df_column)
df_column = df_column.apply(lambda x: tokenizer.tokenize(x))

# 4. Stopwords
# Creating list to store stopwords
stopwords_list = stopwords.words('english')
```

**4: 2-f) Summarizing Review Preprocessing**

In [74]:

```
# Defining the function
def preprocess_review(df_column):
    # 1. Standardizing case
    The previous steps will be gathered into one main function so the above steps can be applied to
    the entire dataset = df_column.str.lower()

In [73]:
```

# Ensuring the relevant packages are imported  
# These were imported individually before but are reminded here if they needed to be imported again  
# df\_column = character encoding(df\_column)  
# from nltk.tokenize import RegexpTokenizer  
# nltk.download('stopwords', quiet=True)  
# from nltk.corpus import stopwords  
# from nltk.stem import WordNetLemmatizer  
 token\_pattern = r"(?u)\b\w+\b"

# Instantiating the tokenizer  
 tokenizer = RegexpTokenizer(token\_pattern)

# Tokenizing  
# df\_column = tokenizer.tokenize(df\_column)  
 df\_column = df\_column.apply(lambda x: tokenizer.tokenize(x))

# 4. Stopwords  
# Creating list to store stopwords  
 stopwords\_list = stopwords.words('english')

# Storing words to add to list of stopwords  
# manual\_stopwords = []
# Adding to list of stopwords
# for word in manual\_stopwords:
# stopwords\_list.append(word)

# Removing stopwords
# df\_column = [token for token in df\_column if token not in stopwords\_list]
df\_column = df\_column.apply(lambda tokens: [token for token in tokens if token not in stopwords\_list])

# 5. Lemmatizing
# Instantiating the Lemmatizer
 lemmatizer = WordNetLemmatizer()

# Lemmatizing words
# df\_column = [lemmatizer.lemmatize(token) for token in df\_column]
df\_column = df\_column.apply(lambda tokens: [lemmatizer.lemmatize(token, pos='v') for token in tokens])

# Returning the preprocessed review
 return df\_column

In [75]:

```
# Verifying the type of X_train['Review']
type(X_train['Review'])
```

Out[75]: pandas.core.series.Series

In [76]:

```
# Testing our function on a few reviews
preprocess_review(X_train['Review'].iloc[95:99])
```

Out[76]:

```
9751
[arc, la, rambla, recommened, hotel, august, 2007, stay, hotel, nights, good, thi
ng, say, clean, safe, room, small, shower, curtain, shower, wash, hold, shower,
hand, guy, reception, helpful, rude, enter, hotel, terrible, smell, follow, go,
stairs, room, balcony, great, air, crap, work, properly, group, room, windows,
air, air, better, room, get, older, room, hotel, decorate, part, hotel, ther, p
erfect, location, la, rambla, shop, im, sure, nicer, hotels, strip]
1213
[wonderful, place, celebrate, 20th, anniversary, place, best, really, quiet, st
aff, perfect, love, ability, walk, snack, day, want, wine, cheese, social, eve
n, want, make, 20th, anniversary, special, casablanca, happen, john, kori]
11506
[style, sophistication, arrive, early, hours, morning, hawaii, tire, jet, lag,
12, hour, flight, soon, perk, saw, hotel, fantastic, property, stay, new, york,
club, floor, room, worth, upgrade, huge, extremely, tastefully, furnish, westi
n, heavenly, bed, say, heavenly, heavenly, bath, westin, promote, little, disap
point, bed, nice, hotel, bathroom, nothing, ordinary, go, say, room, clean, sta
ff, generally, helpful, location, superb, mid, town, time, sqaure, 42nd, stree
t, close especially show shop ground bus tour stop round corner wo
```

```
In [76]: # Testing our function on a few reviews
preprocess_review(X_train['Review'].iloc[95:99])
```

```
Out[76]: 9751
[arc, la, rambla, recommed, hotel, august, 2007, stay, hotel, nights, good, thi
ng, say, clean, safe, room, small, shower, curtain, shower, wash, hold, shower,
hand, guy, reception, helpful, rude, enter, hotel, terrible, smell, follow, go,
stairs, room, balcony, great, air, crap, work, properly, group, room, windows,
air, air, better, room, get, older, room, hotel, decorate, part, hotel, ther, p
erfect, location, la, rambla, shop, im, sure, nicer, hotels, strip]
1213
[wonderful, place, celebrate, 20th, anniversary, place, best, really, quiet, st
aff, perfect, love, ability, walk, snack, day, want, wine, cheese, social, eve
n, want, make, 20th, anniversary, special, casablanca, happen, john, kori]
11506
[style, sophistication, arrive, early, hours, morning, hawaii, tire, jet, lag,
12, hour, flight, soon, perk, saw, hotel, fantastic, property, stay, new, york,
club, floor, room, worth, upgrade, huge, extremely, tastefully, furnish, westi
n, heavenly, bed, say, heavenly, heavenly, bath, westin, promote, little, disap
point, bed, nice, hotel, bathroom, nothing, ordinary, go, say, room, clean, sta
ff, generally, helpful, location, superb, mid, town, time, sqaure, 42nd, stree
t, close, especially, show, shop, greyhound, bus, tour, stop, round, corner, wo
rth, money, overall, great, hotel, stylish, stay, city, cities]
7076
[completely, relax, husband, spend, honeymoon, secrets, week, 14th, ri
de, airport, take, hour, 10, minutes, typical, bus, ride, driver, stop, beers,
way, check, room, ready, eat, lunch, return, desk, room, ready, annemarie, mov
e, room, absolutely, lovely, book, corner, suite, large, balcony, move, floor,
2nd, floor, corner, suite, fabulous, flower, fruit, basket, champagne, sash, do
or, read, honeymooners, bldg, center, resort, convenient, restaurants, activiti
es, room, spotless, daily, maid, service, excellant, room, service, food, good,
take, hour, half, food, restaurants, good, spend, day, pool, beach, palapas, ni
ce, break, sun, careful, mamma, wan, na, potent, stuff, loud, music, pool, min
d, average, monthly, salary, 250, course, ...]
Name: Review, dtype: object
```

```
In [77]: # Applying this to our whole train data
X_train['Review'] = preprocess_review(X_train['Review'])
```

```
In [78]: # Verifying that I correctly have lowercase strings after character encoding, and
is_lowercase(X_train['Review'])
```

The assumption that all reviews are in lowercase is True

```
In [79]: # Now creating a column with of preprocessed reviews without being stored in list
X_train['Review_prep_nolist'] = X_train['Review'].apply(lambda x: ' '.join(x))
```

```
In [80]: # Inspecting the newly created column
X_train[:2]
```

Out[80]:

	Review	Rating	label	original_review	Review_prep_nolist
1282	[gem, pleasantly, surprise, accommodations, helpful, attentive, staff, need, breakfast, good, cookies, room, lovely, clean, comfortable, room, front, bush, street, bite, noisy, away, enjoyment, nice, boutique, hotel, definitely, stay, time]	4	not_detractors	gem, pleasantly surprised accommodations helpful attentive staff need, breakfast good cookies, room lovely clean comfortable room front bush street bite noisy away enjoyment nice boutique hotel definitely stay time,	gem pleasantly surprise accommodations helpful attentive staff need breakfast good cookies room lovely clean comfortable room front bush street bite noisy away enjoyment nice boutique hotel definitely stay time,
	[love, fita, wife, spend, nights, hotel, fita, march, 2008, hotel, staff, fantastic, helpful, pleasant, recommend, hotel, travel,			loved fita wife spent nights hotel fita march 2008. hotel staff fantastic helpful pleasant, recommend hotel traveling amsterdam the hotel	love fita wife spend nights hotel fita march 2008 hotel staff fantastic helpful pleasant recommend hotel

```
In [80]: # Inspecting the newly created column
X_train[:2]
```

Out[80]:

		Review	Rating	label	original_review	Review_prep_nolist
1282		[gem, pleasantly, surprise, accommodations, helpful, attentive, staff, need, breakfast, good, cookies, room, lovely, clean, comfortable, room, front, bush, street, bite, noisy, away, enjoyment, nice, boutique, hotel, definitely, stay, time]	4	not_detractors	gem, pleasantly surprised accommodations helpful attentive staff need, breakfast good cookies, room lovely clean comfortable, room fronting bush street bit noisy did not away enjoyment nice boutique hotel, definitely stay time,	gem pleasantly surprise accomodations helpful attentive staff need breakfast good cookies room lovely clean comfortable room front bush street bite noisy away enjoyment nice boutique hotel definitely stay time
10661		[love, fita, wife, spend, nights, hotel, fita, march, 2008, hotel, staff, fantastic, helpful, pleasant, recommend, hotel, travel, amsterdam, hotel, locate, minute, walk, van, gogh, rijks, museums, 10, minute, walk, leidsplein, variety, store, shop, restaurants, locate, tram, stop, directly, hotel, easy, access, entire, city, walk, fita, way, center, city, 30, minutes, beautiful, accesible, friendly, definately, travel, fita, wonderful, amsterdam]	5	not_detractors	loved fita wife spent nights hotel fita march 2008. hotel staff fantastic helpful pleasant, recommend hotel traveling amsterdam.the hotel located minute walk van gogh rijks museums 10 minute walk leidsplein variety store shop restaurants locate tram stop directly hotel easy access entire city walk fita way center city 30 minutes.beautiful accesible friendly definately travel fita wonderful amsterdam,	love fita wife spend nights hotel fita march 2008 hotel staff fantastic helpful pleasant recommend hotel travel amsterdam hotel locate minute walk van gogh rijks museums 10 minute walk leidsplein variety store shop restaurants locate tram stop directly hotel easy access entire city walk fita way center city 30 minutes beautiful accesible friendly definately travel fita wonderful amsterdam

- Preprocessing test data for later use

```
In [81]: # Creating a duplicate of the review column
X_test['original_review'] = X_test['Review']
```

```
In [82]: # Preprocessing reviews to the test data
X_test['Review'] = preprocess_review(X_test['Review'])
```

```
In [83]: # Now creating a column of preprocessed reviews without being stored in lists, for
X_test['Review_prep_nolist'] = X_test['Review'].apply(lambda x: ' '.join(x))
```

```
In [84]: # Verifying that I correctly have Lowercase strings after character encoding, and
is_lowercase(X_test['Review'])
```

The assumption that all reviews are in lowercase is True

## 4: 2- g) Frequency Distributions

```
In [86]: # Creating an example of Frequency distribution for 1 review
A frequency distribution is a data structure which can be compared to a list displaying how often a specific word or phrase appears.
```

```
Out[86]: FreqDist(doctest) will use the FreqDist() package's 'list' to pass in a single list of words. If early, 'standards', 'I', differ, 'I', 'considerably', 'clean': 1, ...})
```

```
In [87]: # Importing the relevant package for top number of words
from matplotlib.ticker import MaxNLocator
```

```
In [85]: # Importing a function to generate a bar chart for top 10 words
```

```
from nltk import FreqDist
def visualize_top_10(freq_dist, title, rotation):
    # Extracting data for graph
    • FreqDist
        top_10 = list(zip(*freq_dist.most_common(10)))
        tokens = top_10[0]
        counts = top_10[1]

    # Setting up graph and plotting data
    fig, ax = plt.subplots()
    ax.bar(tokens, counts, color='#00000F')
```

In [86]: # Creating an example of Frequency distribution for 1 review  
A frequency distribution is a data structure which can be compared to a list displaying how often a piece of data, or word appears.

Out[86]: FreqDist do expect to use the FreqDist package it always uses a single list of words. It then produces a dictionary-like output of those words and their frequencies {clean': 1, ...}

In [87]: # Importing the relevant package for top number of words  
will visualize the top 10 words to evaluate further what cleaning needs to be done.

In [85]: # Importing relevant packages to visualize top 10 words  
from nltk import FreqDist

```
def visualize_top_10(freq_dist, title, rotation):
    # extracting data for graph
    top_10 = list(zip(*freq_dist.most_common(10)))
    tokens = top_10[0]
    counts = top_10[1]

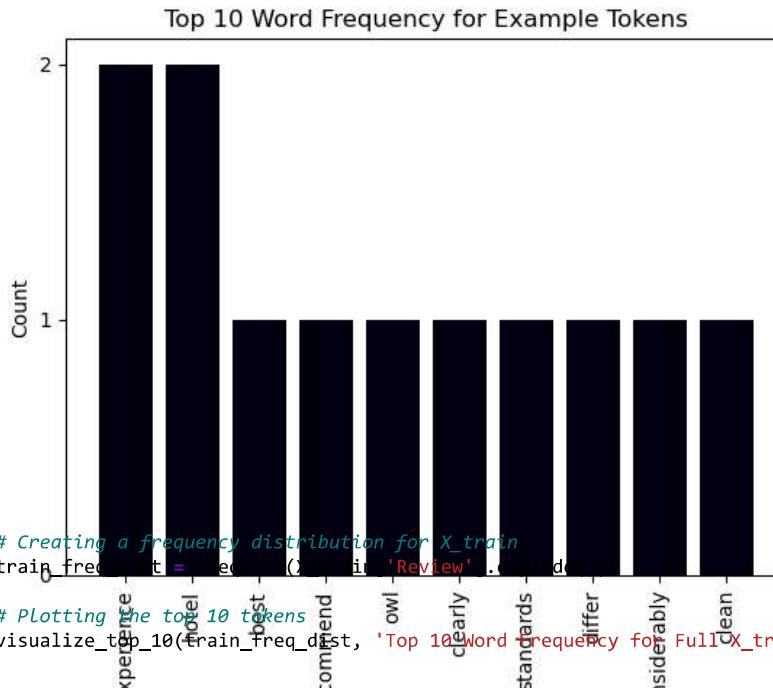
    # Setting up graph and plotting data
    fig, ax = plt.subplots()
    ax.bar(tokens, counts, color="#00000F")
    # ax.set_facecolor('#ffffff')

    # Customizing plot appearance
    ax.set_title(title)
    ax.set_ylabel('Count')

    # Formatting the y-axis Labels to show thousands
    ax.yaxis.set_major_locator(MaxNLocator(integer=True))
    ax.yaxis.set_major_formatter(FuncFormatter(lambda x, pos: '{:,}'.format(int(x))

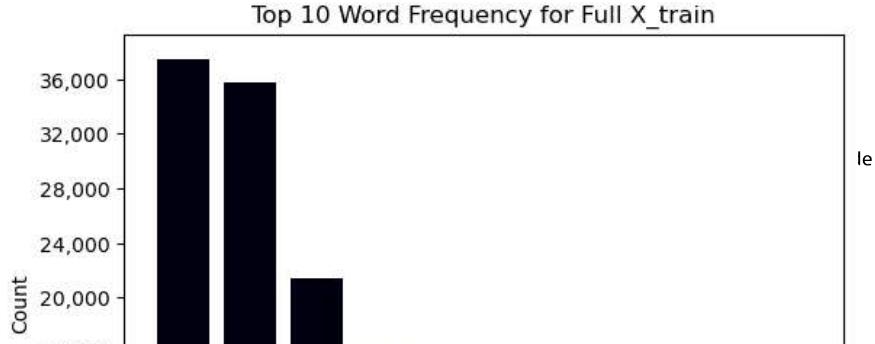
    ax.tick_params(axis='x', rotation=rotation)

visualize_top_10(example_freq_dist, "Top 10 Word Frequency for Example Tokens", 90)
```

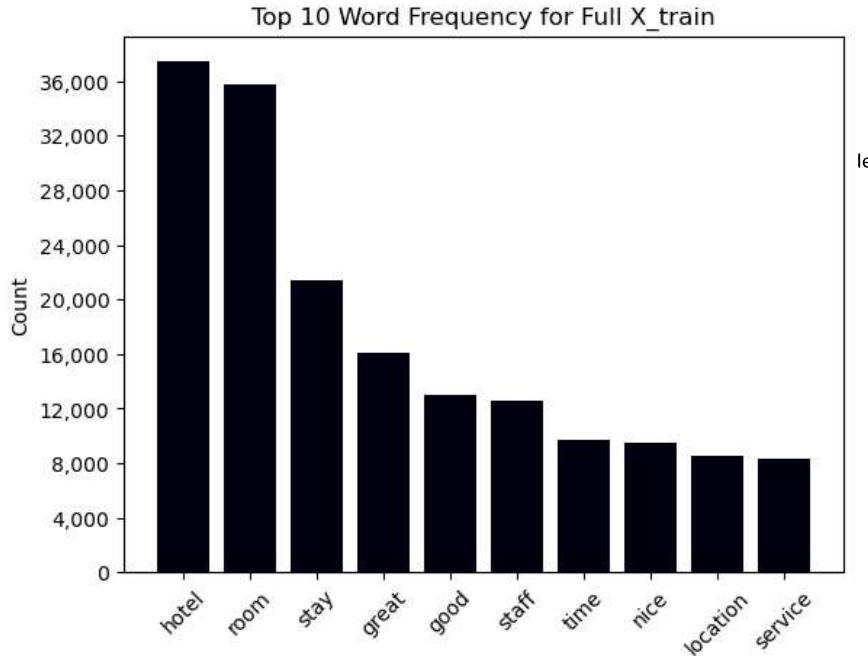


In [88]: # Creating a frequency distribution for X\_train  
train\_freq\_dist = example\_freq\_dist.inplace=True

```
# Plotting the top 10 tokens
visualize_top_10(train_freq_dist, 'Top 10 Word Frequency for Full X_train', 45)
```



```
In [88]: # Creating a frequency distribution for X_train
train_freq_dist = ed.X_train['Review'].value_counts()
# Plotting the top 10 tokens
visualize_top_10(train_freq_dist, 'Top 10 Word frequency for Full X_train', 45)
```



```
In [89]: # Inspecting the most common 20 words
train_freq_dist.most_common(20)
```

```
Out[89]: [('hotel', 37427),
('room', 35770),
('stay', 21391),
('great', 16054),
('good', 13046),
('staff', 12527),
('time', 9687),
('nice', 9459),
('location', 8529),
('service', 8308),
('clean', 8215),
('beach', 7747),
('walk', 7731),
('breakfast', 7539),
('night', 7498),
('day', 7410),
('place', 7377),
('like', 7102),
('food', 7061),
```

```
In [91]: X_train[6750]
```

```
Out[91]: I will also subdivide this by category (detractors/not_detractors) to see if it makes a difference.
```

	Review	Rating	Label	Original_review	Review_prep_nolist
	[gem, pleasantly, surprise, accomodations, helpful, attentive, staff, need, breakfast, good, cookies, room, lovely, clean, comfortable, room, front, bush, street, bite, noisy, away, enjoyment, nice, boutique, hotel, definitely, stay, time]	4	not_detractors	gem, pleasantly surprised accomodations helpful attentive staff need, breakfast good cookies, room lovely clean comfortable, room fronting bush street bite noisy did not away enjoyment nice boutique hotel, definitely stay time,	gem pleasantly surprise accomodations helpful attentive staff need breakfast good cookies room lovely clean comfortable room front bush street bite noisy away enjoyment nice boutique hotel definitely stay time,
1282	[love, fita, wife, spend, nights, hotel, fita, march, 2008, hotel, staff, fantastic, helpful, pleasant, recommend, hotel, travel, amsterdam, hotel, "]			loved fita wife spent nights hotel fita march 2008. hotel staff fantastic helpful pleasant, recommend hotel traveling amsterdam.the hotel located minute walk	love fita wife spend nights hotel fita march 2008 hotel staff fantastic helpful pleasant recommend hotel travel amsterdam hotel

```
In [91]: x[tfidf_start[3] : 6750]
```

Out[91]: I will also subdivide this by category (detractors/not\_detractors) to see if it makes a difference:

		Review	Rating	label	original_review	Review_prep_nolist
		[gem, pleasantly,			surprised	gem pleasantly surprise
		# # Adding in 'label' for filtering			accommodations]	accommodations helpful
		# X_train['label'] == val for val in X_train[val] for val in X			helpful attentive staff	attentive staff need
1282		helpful, attentive, staff, need, breakfast, good, cookies, room, lovely, clean, comfortable, room, front, bush, street, bite, noisy, away, enjoyment, nice, boutique, hotel, definitely, stay, time]	4	not_detractors	need, breakfast good cookies, room lovely clean comfortable, room fronting bush street bit noisy did not away enjoyment nice boutique hotel, definitely stay time,	bite noisy away enjoyment nice boutique hotel definitely stay time,
10661		[love, fita, wife, spend, nights, hotel, fita, march, 2008, hotel, staff, fantastic, helpful, pleasant, recommend, hotel, travel, amsterdam, hotel, locate, minute, walk, van, gogh, rijks, museums, 10, minute, walk, leidsplein, variety, store, shop, restaurants, locate, tram, stop, directly, hotel, easy, access, entire, city, walk, fita, way, center, city, 30, minutes, beautiful, accesible, friendly, definately, travel, fita, wonderful, amsterdam]	5	not_detractors	loved fita wife spent nights hotel fita march 2008. hotel staff fantastic helpful pleasant, recommend hotel traveling amsterdam.the hotel located minute walk van gogh rijks museums 10 minute walk leidsplein variety stores shops restaurants located.there tram stop directly hotel, easy access entire city, did walk fita way center city 30 minutes.beautiful accesible friendly definately travel fita wonderful amsterdam,	love fita wife spend nights hotel fita march 2008 hotel staff fantastic helpful pleasant recommend hotel travel amsterdam hotel locate minute walk van gogh rijks museums 10 minute walk leidsplein variety store shop restaurants locate tram stop directly hotel easy access entire city walk fita way center city 30 minutes beautiful accesible friendly definately travel fita wonderful amsterdam
17908		[great, modern, hotel, fantastic, price, stay, night, begin, april, alot, review, say, bad, area, really, literally, 100m, away, nice, redevelop, area, 10min, walk, beach, partner, arrive, late, barcelona, decide, walk, hotel, bus, station, order, bar, quite, long, walk, 50mins, roughly, really, great, thing, order, sight, tourists, dont, usually, arrive, check, quick, painless, room, really, really, nice, fantastic, modern, design, little, extras, want, order, barcelona, partner, opt, bus, turistic, stop, far, hotel, go, virtually, worth, short, visit, wish, spend, longer, hotel, wish, swim, pool, roof, open, overall, fantastic, hotel, price, pay, room, dont, distance, centre, transport, options]	5	not_detractors	great modern hotel fantastic price stayed just night beginning april, alot reviews say bad area really didn't, literally 100m away nice redeveloped area 10min walk beach.me partner arrived late barcelona decided walk hotel bus station order barings, quite long walk 50mins roughly really great thing order sights tourists dont usually, arrived check quick painless room really really nice, fantastic modern design little extras want, order barcelona partner opted bus turistic, stop not far hotel goes virtually, worth short visit,i wish spent longer hotel, wish swimming pool roof open, overall fantastic hotel price paid room dont distance centre transport options	great modern hotel fantastic price stay night begin april alot review say bad area really literally 100m away nice redevelop area 10min walk beach partner arrive late barcelona decide walk hotel bus station order bar quite long walk 50mins roughly really great thing order sight tourists dont usually arrive check quick painless room really really nice fantastic modern design little extras want, order barcelona partner opt bus turistic stop far hotel go virtually worth short visit wish spend longer hotel wish swim pool roof open overall fantastic hotel price pay room dont distance centre transport options

```
In [93]: # Defining function to plot 2 visualizations
```

```
# Creating two columns
def two_subplots():
```

```
# Reminding plt figure(figsize=(15, 9))
# fig=plt.figure(figsize=(15, 9))
# fig.set_tight_layout(True)
# fig.set_facecolor('#254858')
# custom_set_color(['#60A787', '#DEB887']) #green, gold
custom_set_color(['#254858', '#AA4255']) # green, dark red

# custom=colored_subplots(gs[0, 0], "#00A887", 0, col_0
ax2 = fig.add_subplot(gs[0, 1]) #row 0, col 1
return fig, [ax1, ax2]
```

```
# Plotting the graph
```

```
def plot_distribution_by_sentiment(X_version, column, axes, rotation):
    for index, category in enumerate(X_version['label'].unique()):
        # Calculating frequency distribution for this subset
        all_words = X_version[X_version['label'] == category][column].explode()
        freq_dist = FreqDist(all_words)
        top_10 = list(zip(*freq_dist.most_common(10)))
        tokens = top_10[0]
```

```
In [93]: # Defining function to plot 2 visualizations
# Creating two columns
def two_subplots():

In [92]: # Remaking custom colors
# Defining custom layout (True)
# custom set facecolor '#2F4858' '#DEB887' ] #green, gold
custom_color = ['#2F4858', '#DEB887'] #green, dark red
custom_color += [ '#00A425', '#AA4255' ] # green, dark red

# custom fig.add_subplot(gs[0,0]) #row 0, col 0
ax1 = fig.add_subplot(gs[0,0]) #row 0, col 0
ax2 = fig.add_subplot(gs[0,1]) #row 0, col 1
return fig, [ax1, ax2]

# Plotting the graph
def plot_distribution_by_sentiment(X_version, column, axes, rotation):
    for index, category in enumerate(X_version['label'].unique()):
        # Calculating frequency distribution for this subset
        all_words = X_version[X_version['label'] == category][column].explode()
        freq_dist = FreqDist(all_words)
        top_10 = list(zip(*freq_dist.most_common(10)))
        tokens = top_10[0]
        counts = top_10[1]

        # Setting up a plot
        ax = axes[index]
        ax.bar(tokens, counts, color=custom_color[index])

        # Setting background color
        ax.set_facecolor('#2F4858')

        # Customizing plot appearance
        title = "Word Frequency for:" + category
        ax.set_title(f'{title}', fontsize=17, color='white')
        ax.set_ylabel("Count", fontsize=12, color='white')

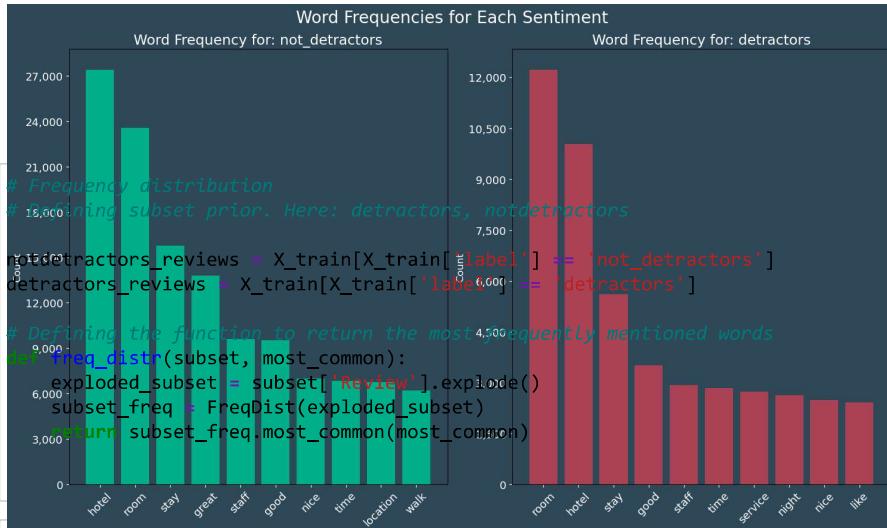
        # Setting tick color
        ax.tick_params(axis='x', colors='white', labelsize=13)
        ax.tick_params(axis='y', colors='white', labelsize=13)

        # Formatting the y-axis Labels to show thousands
        ax.yaxis.set_major_locator(MaxNLocator(integer=True))
        ax.yaxis.set_major_formatter(FuncFormatter(lambda x, pos: '{:,}'.format(x)))
        ax.tick_params(axis='x', rotation=rotation)

fig, axes = two_subplots()
plot_distribution_by_sentiment(X_train, 'Review', axes, 45)
fig.suptitle('Word Frequencies for Each Sentiment', fontsize=20, color='white')

# Saving the plot as a PNG with a transparent background
plt.savefig('images/word_freq.png', transparent=True)

plt.show()
```



```
In [94]: # Frequency distribution
# training subset prior. Here: detractors, notdetractors
notdetractors_reviews = X_train[X_train['label'] == 'not_detractors']
detractors_reviews = X_train[X_train['label'] == 'detractors']

# Defining the function to return the most frequently mentioned words
def freq_distr(subset, most_common):
    exploded_subset = subset['Review'].explode()
    subset_freq = FreqDist(exploded_subset)
    return subset_freq.most_common(most_common)

In [95]: # Getting frequency distribution for top 20 strings of detractors
freq_distr(detractors_reviews, 20)
```

```
In [94]: # Frequency distribution
# Beginning subset prior. Here: detractors, notdetractors
notdetractors_reviews = X_train[X_train['label'] == 'not_detractors']
detractors_reviews = X_train[X_train['label'] == 'detractors']

# Defining the function to return the most frequently mentioned words
def freq_distr(subset, most_common):
    exploded_subset = subset['Review'].explode()
    subset_freq = FreqDist(exploded_subset)
    return subset_freq.most_common(most_common)
```

```
In [95]: # Getting frequency distribution for top 20 strings of detractors
freq_distr(detractors_reviews, 20)
```

```
Out[95]: [('room', 12221),
('hotel', 10028),
('stay', 5608),
('good', 3507),
('staff', 2920),
('time', 2840),
('service', 2722),
('night', 2630),
('nice', 2479),
('like', 2406),
('resort', 2362),
('day', 2335),
('get', 2330),
('go', 2310),
('food', 2286),
('beach', 2286),
('great', 2272),
('place', 2081),
('clean', 2068),
('say', 2046)]
```

Now reviewing the first 20 words for reviews coming from detractors help us confirm that the first top 5 words are as represented in both categories. They are most common words to refer to a hotel stay. They should be removed for clearer analysis.

```
In [96]: # # Applying this to our whole dataset
# X_train['Review'] = X_train['Review'].apply(lambda x: preprocess_review(x))
```

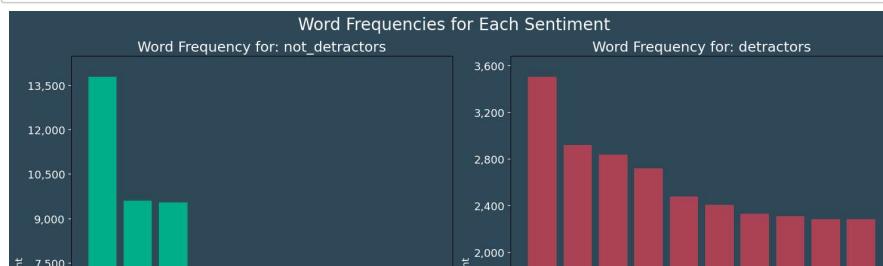
```
In [97]: # Including these new words to a new list of stopwords
new_stopwords = ['hotel', 'room', 'night', 'day', 'stay', 'resort', 'place']
```

```
In [98]: # Calling the previously defined function to remove stopwords
X_train['Review'] = X_train['Review'].apply(lambda x: remove_stopwords(x, new_sto
```

```
In [99]: # Calling the function again to review the new split of top 10 words for each sentiment
fig, axes = two_subplots()
plot_distribution_by_sentiment(X_train, 'Review', axes, 45)
fig.suptitle('Word Frequencies for Each Sentiment', fontsize=20, color='white')

# Saving the plot as a PNG with a transparent background
plt.savefig('images/word_freq.png')

plt.show()
```

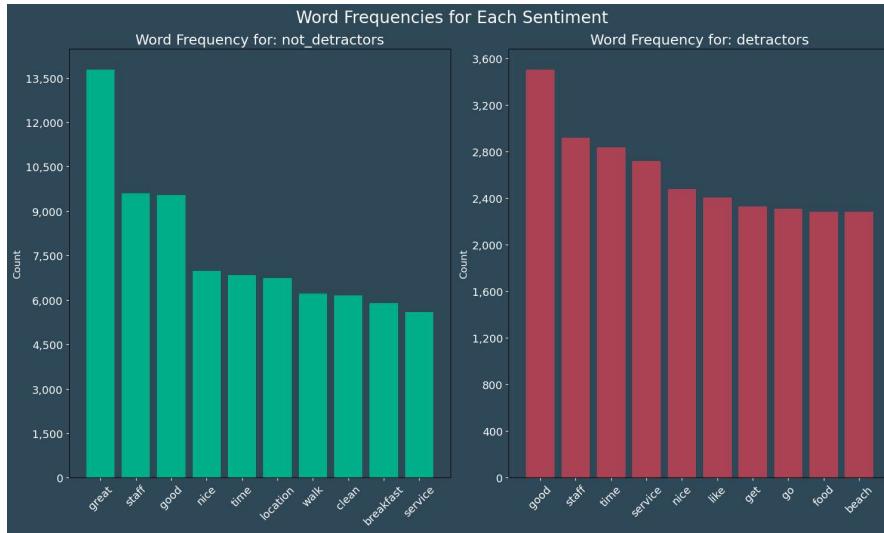


```
In [99]: # Calling the function again to review the new split of top 10 words for each sentiment

fig, axes = two_subplots()
plot_distribution_by_sentiment(X_train, 'Review', axes, 45)
fig.suptitle('Word Frequencies for Each Sentiment', fontsize=20, color='white')

# Saving the plot as a PNG with a transparent background
plt.savefig('images/word_freq.png')

plt.show()
```



I will filter on reviews that include specific words such as 'time' to review if they should be considered stopwords.

```
In [100]: # Defining the word to search
word_search = 'time'

# Printing the results
X_train[X_train['Review_prep_nolist'].str.contains(word_search)][:3]

# Another way of searching directly in the column 'Review' in case several words
# word_researched = 'time'
# filtered_on_word = X_train[X_train['Review'].apply(lambda review_list: any(wor
# filtered_on_word
```

[hat, family, husband, children, disable, june, 23, 30, 2007, ok, time, hot, water, time, contact, service, desk, advise, respond, attitude, off, end, happy, reflect, tremendously, service, food, terrible, sick, end, hospital, bar graph, beg, doctor, discharge, early, children, claim, sick, Word clouds visually represent the frequency of words in a given text, with more frequent words occurring displayed in larger font size. This allows for quick intuition about what is important in the text. husband 3 children disabled, stayed resort june 23-30 2007. room ok times did not hot water, time contacted service desk advise respond attitude, staff working resort looked not happy reflected tremendously service, food terrible, sick ended hospital beg doctor discharge early children, hotel claimed no sick imagination, rep suggested eat suggestion, hotel aware hospital stay ended pok hotel doctor office equipped handle bacterial infection, no hotel contacted doing hat resort family husband children disable stay resort june 23 30 2007 room ok time hot water time contact service desk advise respond attitude staff work resort look happy reflect tremendously service food terrible sick end hospital beg doctor discharge early imagination rep suggest eat bread week nice suggestion hotel aware hospital stay ended pok hotel doctor office equipped handle bacterial infection, no hotel contacted doing pick hotel doctor office equip handle bacterial infection hotel contact

```
In [101]: # Installing wordcloud
# Time is so commonly used to describe either category; such as 'next time' to describe a 'times
# !pip install --trusted-host pypi.org --trusted-host pypi.python.org --trusted-
```

understand through visualizations and bigrams.

```
In [102]: # Importing relevant packages
        from matplotlib.colors import LinearSegmentedColormap

        # Defining a colormap that interpolates between the two colors
        custom_colors_detractors = ['#BFA5A7', '#AA4255']
        # custom_colors_detractors = ['#FEEBC4', '#AA4255']
```



```
In [105]: # Examining the 'Review_prep_holistic' column
X_train.columns
```

```
Out[105]: Index(['Review', 'Rating', 'label', 'original_review', 'Review_prep_holistic'], d
type='object')
```

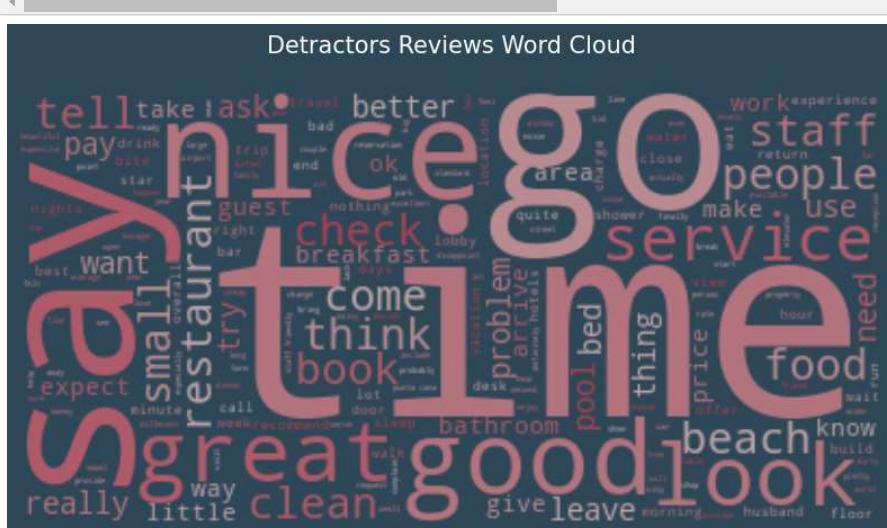
```
In [106]: # Impacting relevance
ask available
# Concatenating detractors and not_detractors reviews separately
detractors_reviews = ' '.join(X_train[X_train['label'] == 'detractors'][['Review']]
not_detractors_reviews = ' '.join(X_train[X_train['label'] == 'not_detractors'][['Review']])
```

Place, resort and time are really stand out among some of the most talked about topics in all reviews. However, it is difficult to distinguish whether these reviews indicate anything mentioned by detractors or not so I will divide them into 2 categories:

```
wordcloud = WordCloud(width=400, height=200, background_color=background_col
```

```
# Displaying the generated word cloud using matplotlib
fig = plt.figure(figsize=(10, 5))
fig.set_facecolor('#2F4858')
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.title(title, color='white', fontsize=15, pad=25)
plt.show()
```

```
# Plotting detractors reviews word cloud
wordcloud_graph(detractors_reviews, 'Detractors Reviews Word Cloud', cmap_detract
```



```
In [107]: # Plotting not_detractors reviews word cloud  
wordcloud_graph(not_detractors_reviews, 'Not detractors Reviews Word Cloud', cmap
```



```
In [107]: # Plotting not_detractors reviews word cloud  
wordcloud_graph(not_detractors_reviews, 'Not detractors Reviews Word Cloud', cmap
```



While the wordcloud for not\_detractors seem to standout more, the ones from detractors look like they are a bit harder to distinguish due to the number of times "time" is mentioned. I will try removing it from this extract, along with say and go to better understand the themes.

```
In [108]: # Making a copy of the list of words
detractors reviews modified = detractors reviews
```

```
In [109]: # Defining the strings to remove  
strings to remove = ['time', 'say', 'go', 'nice', 'great', 'od', 'fo']
```

```
In [110]: # Removing them
for string in strings_to_remove:
    detractors_reviews_modified = detractors_reviews_modified.replace(string, '')
# detractors reviews modified
```

```
In [111]: # Plotting detractors reviews word cloud  
wordcloud_graph(detractors_reviews_modified, 'Detractors Reviews Word Cloud', cm=  
plt.savefig('images/wordcloud_detractors.png')
```



```
In [112]: # Some words were researched individually as well to ensure  
X_train[X_train['Review prep no.list'].str.contains('service')][::1]
```

Out[112]: 

Review	Rating	label	original_review	Review_prep_nolist
[hat, family, husband, children, disable, june, 23, 30, 2007, ok time hot water, time, contact, service, desk, advise, respond, attitude, staff, work, look, happy, reflect, tremendously, service, food, terrible, sick, end, hospital, beg, doctor, discharge, early, children, claim, sick, imagination, rep, suggest, eat, bread,	5	hated-resort-family-5-husband-3-children-disabled-stayed-resort-june-23-30-2007-room-ok-times-did-not-hot-water-time-contact-service-desk-advice-respond-attitude-staff-working-resort-looked-not-happy-reflected-tremendously-service-food-terrible-sick-ended-hospital-beg-doctor-discharge-early-children-hotel-claimed-no-sick-imagination-rep-suggested-eat-bread-week-nice-suggestion-hotel-aware-hospital-stay-ambulance-pick-hotel-doctor-office-not-equipped-hands-bacterial-infection	hated resort family 5. husband 3 children disabled, stayed resort june 23-30 2007. room ok times did not hot water, time contacted service desk advise respond attitude, staff working resort looked not happy reflected tremendously service, food terrible, sick ended hospital beg doctor discharge early children, hotel claimed no sick imagination, rep suggested eat bread week nice suggestion, hotel aware hospital stay ambulance pick hotel doctor office not equipped hands bacterial infection	hat resort family husband children disable stay resort june 23 30 2007 room ok time hot water time contact service desk advise respond attitude staff work resort look happy reflect tremendously service food terrible sick end hospital beg doctor discharge early children hotel claim sick imagination rep suggest eat bread week nice suggestion hotel aware hospital stay ambulance pick hotel doctor

In [112]:

	Review	Rating	label	original_review	Review_prep_nolist
6730	hat, family, husband, children, disable, june, 23, 30, 2007, time, not water, time, contact, service, desk, advise, respond, attitude, staff, work, look, happy, reflect, tremendously, service, food, terrible, sick, end, hospital, beg, doctor, discharge, early, children, claim, sick, imagination, rep, suggest, eat, bread, week, nice, suggestion, aware, hospital, ambulance, pick, doctor, office, equip, handle, bacterial, infection, contact, actual, fact, treat, garbage, kid, club, kid, club, children, active, love, swim, play, sport, interact, children, want, kid, club, nothing, kid, kid, sing, songs, circle, break, heart, children, fun, years, away, children, travel, family, cheap, especially, take, years, put, away, money, little, ...]	1	detractors	hated resort family. - husband 3 children disabled, stayed resort june 23-30 2007. room ok times did not hot water, time contacted service desk advise respond attitude, staff working resort looked not happy reflected tremendously service, food terrible, sick ended hospital beg doctor discharge early children, hotel claimed no sick imagination, rep suggested eat bread week nice suggestion, hotel aware hospital stay ambulance pick hotel doctor office not equipped handle bacterial infection, no hotel contacted doing actual fact treated garbage, kids club, kids club, children active love swim play sports interact children did not want kids club, did nothing, kids night, kids, night sang songs circle, broke heart children not having fun, years away children, having travel family not cheap, especially takes years putting away money little little order enjoy week vacation, disappointed hotel not recommend unless young going bunch girls guys, like types travelers, lived dr travelled dr 15 times treated poorly did, ignored posting site warning figured people picky, realize hotels dr just not worth,	hat resort family husband children disable stay resort june 23 30 2007 room ok time hot water contact service desk advise respond attitude staff work resort look happy reflect tremendously service food terrible sick end hospital beg doctor discharge early children hotel claim sick imagination rep suggest eat bread week nice suggestion hotel aware hospital stay ambulance pick hotel doctor office equip handle bacterial infection hotel contact actual fact treat garbage kid club kid club children active love swim play sport interact children want kid club nothing kid night kid night sing songs circle break heart children fun years away children travel family cheap especially take years put away money little little order enjoy week vacation disappoint hotel recommend unless young go bunch girls guy like type travelers live dr travel dr 15 time treat poorly ignore post site warn figure people picky realize hotels dr worth

Themes that come up the most for the detractors' reviews are:

- the overall service and staff friendliness:
  - service
  - staff
  - people
- the view from the rooms
- the decor
  - look
- cleanliness
  - clean
- problem solving
  - problem

Coming up with these themes required extra research though through individual words, so I will try to use bigrams to make the themes more obvious.

Other smaller negative items:

## 4.2 Bigrams

- Bathroom

The mention of 'pay' is indicated. When contrasted with the 'value' mentioned in reviews coming from non-detectors, this indicates guests are okay with paying higher rates as long as they feel like it represents a good value for money.

In [113]:

# Importing relevant package

As initial recommendations, import `nltk` to reduce the number of reviews from detractors:

1. Provide training sessions to enhance the friendliness of the staff
2. Re-define cleanliness standards
3. Revise the aesthetics of the lobby and rooms—appearance is crucial

4. Cultivate a problem-solving environment for staff

I will define a function to review bigrams and be able to filter in the future by the category I need.

Coming up with these themes required extra research through individual words, so I will try to use bigrams to make the themes more obvious.

Other smaller negative items:

#### 4.2 Bigrams

- Bathroom

The mention of 'pay' is indicated. When contrasted with the 'value' mentioned in reviews coming from non-detectors, this indicates guests are okay with paying higher rates as long as they feel like it represents a good value for money.

**Air Reviews**

```
In [113]: # Importing relevant package
As initial recommendations, I plan to reduce the number of reviews from detractors:
1. Provide training sessions to enhance the friendliness of the staff
# Storing n-locations BigramAssocMeasures into variable
bigram_measures = nltk.collocations.BigramAssocMeasures()
3. Revise the aesthetics of the lobby and rooms—appearance is crucial
4. Cultivate a problem-solving environment for staff
I will define a function to review bigrams and be able to filter in the future by the category I need.
```

```
In [114]: # Defining a function to review bigrams
def bigram_review(text, top_n):
    # Creating a finder and passing it the words of reviews summarized as 1 list
    text_finder = BigramCollocationFinder.from_words(text.sum())
    text_scored = text_finder.score_ngrams(bigram_measures.raw_freq)
    return text_scored[:top_n]
```

I am starting with the top 20 bigrams overall.

```
In [115]: # Calling the function to review the top 20 bigrams overall
bigram_review(X_train['Review'], 20)
```

```
Out[115]: [('great', 'location'), 0.001195786887061),
  ('staff', 'friendly'), 0.0011807859762914177),
  ('punta', 'cana'), 0.0008879110517424272),
  ('walk', 'distance'), 0.000804334548883325),
  ('friendly', 'helpful'), 0.0007857619926924135),
  ('highly', 'recommend'), 0.0007250440205298179),
  ('staff', 'helpful'), 0.0005964647853619684),
  ('minute', 'walk'), 0.0005750349128339935),
  ('read', 'review'), 0.000535032484115107),
  ('location', 'great'), 0.0005036020044074104),
  ('air', 'condition'), 0.0005000303589860813),
  ('make', 'sure'), 0.0004978873717332838),
  ('san', 'juan'), 0.0004943157263119546),
  ('great', 'time'), 0.0004857437773007647),
  ('good', 'value'), 0.00042859745055949826),
  ('good', 'location'), 0.0004228828178853716),
  ('breakfast', 'buffet'), 0.00038716636367208007),
  ('food', 'good'), 0.0003850233764192826),
  ('new', 'york'), 0.000379308743745156),
  ('value', 'money'), 0.0003757370983238268)]
```

```
In [117]: # Calling the bigram function on detractors
This started providing better insights in addition to what was identified for review, it can be
detected_providing_better_insights_in_review_to_what_was_identified_for_review, it can be
identified here that what is most appreciated by guests are:
# Calling the bigram function on not detractors
not_detractor_bigrams = bigram_review(X_train_notdetractors['Review'], 100)
```

- Friendly and helpful staff

```
In [118]: # Creating wordclouds for bigrams in WordClouds
def wordcloud_bigrams(bigram_list, title, colormap, background_color):
```

# Converting bigrams to a dictionary for WordCloud

```
In [116]: # WordCloudData is {bigram: freq for bigram in bigram_list}
X_train_detectors = X_train[X_train['label'] == 'detectors']
X_train_notdetectors = X_train[X_train['label'] != 'detectors']
wordcloud = WordCloud(width=800, height=400, background_color=background_color,
```

# Plotting WordCloud

```
plt.figure(figsize=(10,5))
plt.imshow(wordcloud, interpolation='bilinear')
plt.title(title, fontsize=15, pad=20)
plt.axis('off')
plt.show()
```

```
In [117]: # Calling the bigram function on detractors
This started providing better insights on what was identified for detractors, it can be
detectors[bigrams_topkey].bigram_review(X_train_detractors['Review'], 100)
identified here that what is most appreciated by guests are:
# Calling the bigram function on not detractors
not_detectors[bigrams_topkey].bigram_review(X_train_notdetractors['Review'], 100)
    • Friendly and helpful staff
```

```
In [118]: # Creating value for money would be felt  
def wordcloud_bigrams(bigram_list, title, colormap, background_color)
```

```
In [116]: # Converting bigrams to a dictionary for WordCloud  
# Credit: loadsdata, bigramforfreq, freq_topnfromviewbigram.list
```

```
In [110]: # WordClouds & Diagrams
X_train_detractors = X_train[X_train['label'] == 'detractors']
X_train_no_detractors = X_train[X_train['label'] != 'detractors']

wordcloud = WordCloud(width=800, height=400, background_color=background_color)

# Plotting WordCloud
plt.figure(figsize=(10,5))
plt.imshow(wordcloud, interpolation='bilinear')
plt.title(title, fontsize=15, pad=20)
plt.axis('off')
plt.show()
```

```
In [119]: # Calling it on the detractors_bigrams_100 List  
wordcloud_bigrams(detractors_bigrams_100, 'Detractors Reviews: Bigrams', cmap_det  
  
# Calling it on the notdetractors_bigrams_100 List  
wordcloud_bigrams(notdetractors_bigrams_100, 'Not Detractors Reviews: Bigrams', c
```

## Detractors Reviews: Bigrams



## Not Detractors Reviews: Bigrams

```
# Calling the bigram function on not detractors
notdetractors.bigrams(500) -> higram_review(X_train.notdetractors[["Review"]], 500)
```

```
In [121]: # Call the function on the detector bigrams list
```

```
# Get them in the notdetractor bigrams_500 list
wordcloud_bigrams(nodetractors_bigrams_500, 'Not Detractors_Reviews: Bigrams',
```

## Detractors Reviews: Bigrams



```
In [120]: # Calling the bigram function on detractors
detractors_bigrams_500 = bigram_reviewer.print_top_n(detractors,
                                         n=500)

# Calling the bigram function on not
notdetractors_bigrams_500 = bigram_reviewer.print_top_n(notdetractors,
                                         n=500)

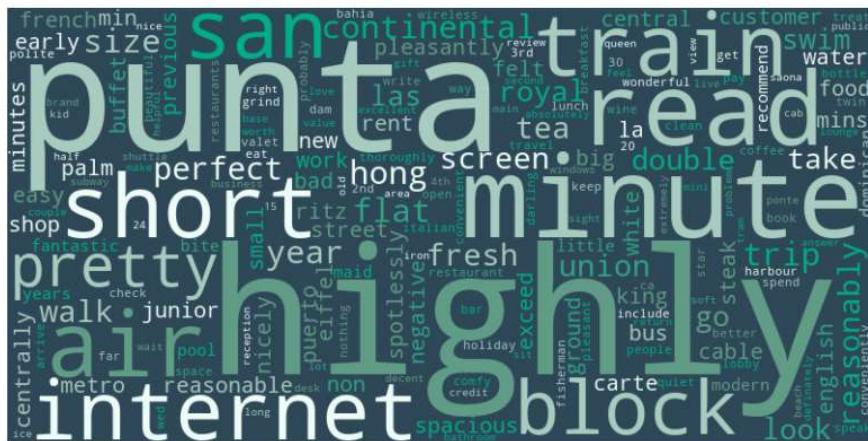
In [121]: # Calling up on the detractor bigrams
wordcloud_bigrams(detractors_bigrams_
                   .head(10),
                   title='Detractor Bigrams')

# Calling up on the notdetractor bigrams
wordcloud_bigrams(notdetractors_bigrams_
                   .head(10),
                   title='Notdetractor Bigrams')
```

## Detractors Reviews: Bigrams



## Not detractors Reviews: Bigrams



There are still some similarities among detractors and not detractors reviews referring to Punta Cana, and walking. Let's research these specific mention to get a better understanding. Out of curiosity, I drew diagrams for Punta Cana specifically to get an understanding of why they are

```
In [122]: # Research walking and punta in detractors bigrams 500
```

```
In [124]: # Calling the bigram function on punta
```

```
def research_bigram(word, bigram_list):
    word_researched = [bigram for bigram in bigram_list if word in bigram[0]]
```

# Caret return the bigram function on not detractors

```
notdetractors_bigrams_punta = bigram_review(X_train_notdet)
```

print(pocoranch.bigrnm['Inunt']) dotplot

```
# print(research_programs, prefix = "defactors_programs_")
```

```
In [125]: # calling it on the detectors test
```

```
wordcloud bigrams(detractors bigrams punta, 'Punta Cana De
```

`l = [100] * 1000000000`

# Calling it on the `notpositive_bigrams_punta` list

```
wordcloud(bigrams[notdefactors], bigrams.punta, 'Punta Cana Not  
#paint(az) bigram hidcam('wulbibat notañitius bigram 500))
```



Cana, and walking. Let's research these specific mention to get a better understanding. Out of curiosity, I drew biograms for Punta Cana specifically to get an understanding of why they are

```
In [122]: # Research walking and punta in detractors_bigrams_500
In [124]: # Calling the bigram function on punta
detractors_bigrams_500,=bigram_review(X_train_detractors[X_train_detractors['Punta Cana Not Detractors Reviews'].str.contains('punta')], word_researched = [bigram for bigram in bigram_list if word in bigram[0]])
# Call the function to research punta on not detractors
notdetractors_bigrams_punta = bigram_review(X_train_notdetractors[X_train_notdetractors['Punta Cana Not Detractors Reviews'].str.contains('punta')], word_researched = [bigram for bigram in bigram_list if word in bigram[0]])
# print(research_bigram('punta', detractors_bigrams_500))
In [125]: #((('punta'))t'earths,det00099021799216478])
wordcloud_bigrams(detractors_bigrams_punta, 'Punta Cana Detractors Reviews: Bigram')
In [123]: # notdetractors_bigrams_500
# Calling it on the notpositive_bigrams_punta List
wordcloud_bigrams(notdetractors_bigrams_punta, 'Punta Cana Not Detractors Reviews')
# print(research_bigram('walking', notpositive_bigrams_500))
print(research_bigram('punta', notdetractors_bigrams_500))
Punta Cana Detractors Reviews: Bigram
```



## Punta Cana Not Detractors Reviews: Bigrams



```
In [127]: # Combining bigram words into a single string for each sentiment
detractor_bigram_20 = [''.join(bigram[0]) for bigram in detractors_bigrams_20]
It looks like Punta Cana is one of the top visited resorts or one of those generating the highest
notdetractor_bigram_20 = [''.join(bigram[0]) for bigram in notdetractors_bigrams_
number of reviews which explains why the name comes up so often.
```

# Reviewing what the newly created Lists Look Like  
By separating only for Punta Cana the negatives and not negative words, I can get better insights:  
#dector diagram 20  
horseback riding is one of the top activities, guests appreciate the all inclusive service, and the hotel

In [128]: # Plotting the detractor bigrams as a bar chart

I will now try to visualize the top 20 positive and negative bigrams overall to close on this topic.

# Extracting words from the document. We will use positive and negative bigrams overall to close on this topic.

```
[100]: bigrams, values = zip(*detractors_bigrams_20)
```

```
In [126]: # Storing the top 20 bigrams for each sentiment  
# fest+is+terrible and val+helps+in+decreasing+neg+med+replaced+on+values
```

for i in range(len(values) - 1, 0, -1):  
 if values[i] < values[i - 1]:  
 values[i], values[i - 1] = values[i - 1], values[i]

```
sorted_bigrams = [bigrams[i] for i in sorted_indices]
```

```
# Create a bar chart
```

```
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#2F4858')
```

```
| ax.set_facecolor('#2F4858')
```

host:8888/notebooks/tripadvisor\_reviews\_detractors.ipynb

```
In [127]: # Combining bigram words into a single string for each sentiment
detractor_bigram_20 = [' '.join(bigram[0]) for bigram in detractors_bigrams_20]
It looks like Punta Cana is one of the top visited resorts or one of those generating the highest
notdetractor_bigram_20 = [' '.join(bigram[0]) for bigram in notdetractors_bigrams_
number of reviews which explains why the name comes up so often.

# Reviewing what the newly created Lists Look Like
By separating only for Punta Cana the negatives and not negative words, I can get better insights:
# detractor_bigram_20
horseback riding is one of the top activities, guests appreciate the all inclusive service, and the hotel
quests are referring to is the Bahia Principe hotel and resort.

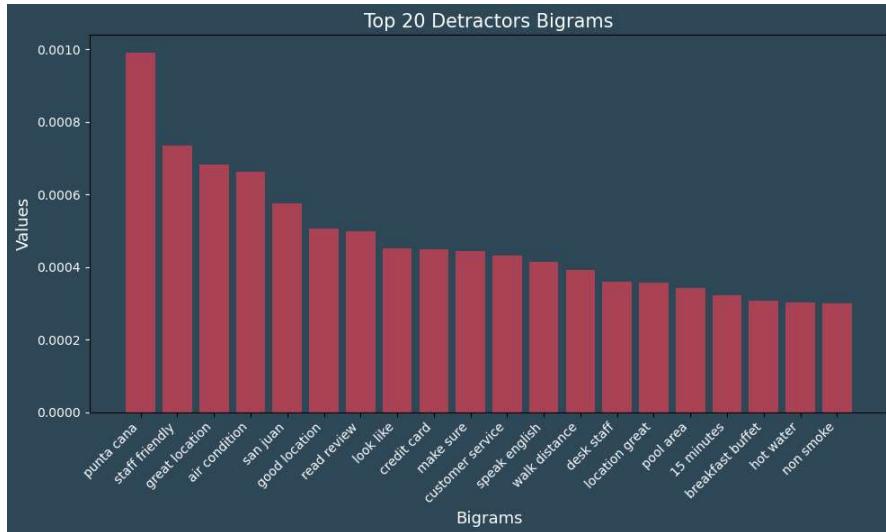
In [128]: # Plotting the detractor bigrams as a bar chart

I will now try to visualize the top 20 positive and negative bigrams overall to close on this topic.

# Extracting bigrams and values
bigrams, values = zip(*detractors_bigrams_20)

In [126]: # Storing the top 20 bigrams for each sentiment
detractors_bigrams_values = detractors_bigrams_20[sorted on values]
sorted_indices = sorted(range(len(values)), key=lambda k: values[k], reverse=True)
sorted_bigrams = [bigrams[i] for i in sorted_indices]
sorted_values = [values[i] for i in sorted_indices]

# Create a bar chart
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')
ax.bar(range(len(sorted_values)), sorted_values, align='center', color='AA4255')
ax.set_xticks(range(len(sorted_values)))
ax.set_xticklabels([f'{bigram[0]} {bigram[1]}' for bigram in sorted_bigrams], rotation=90)
ax.set_xlabel('Bigrams', color='white', fontsize=13)
ax.set_ylabel('Values', color='white', fontsize=13)
ax.tick_params(axis='y', colors='white')
ax.set_title('Top 20 Detractors Bigrams', color='white', fontsize=15)
plt.tight_layout()
plt.show()
```



```
In [129]: # Plotting the not_detractor bigrams as a bar chart

# Extracting bigrams and values
bigrams, values = zip(*notdetractors_bigrams_20)

# Sort bigrams and values in descending order based on values
sorted_indices = sorted(range(len(values)), key=lambda k: values[k], reverse=True)
sorted_bigrams = [bigrams[i] for i in sorted_indices]
sorted_values = [values[i] for i in sorted_indices]

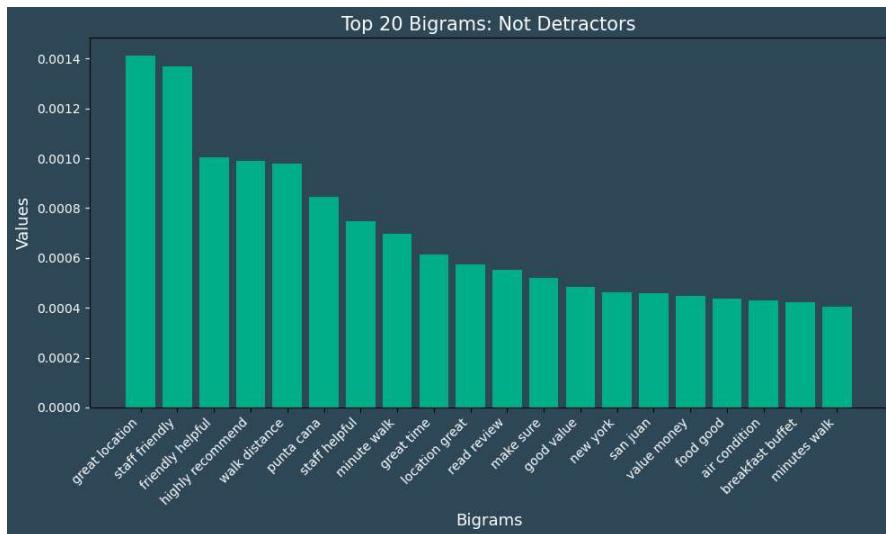
# Create a bar chart
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')
ax.bar(range(len(sorted_values)), sorted_values, align='center', color='00AF87')
ax.set_xticks(range(len(sorted_values)))
ax.set_xticklabels([f'{bigram[0]} {bigram[1]}' for bigram in sorted_bigrams], rotation=90)
ax.set_xlabel('Bigrams', color='white', fontsize=13)
ax.set_ylabel('Values', color='white', fontsize=13)
ax.tick_params(axis='y', colors='white')
ax.set_title('Top 20 Bigrams: Not Detractors', color='white', fontsize=15)
```

```
In [129]: # Plotting the not_detractor bigrams as a bar chart

# Extracting bigrams and values
bigrams, values = zip(*notdetractors_bigrams_20)

# Sort bigrams and values in descending order based on values
sorted_indices = sorted(range(len(values)), key=lambda k: values[k], reverse=True)
sorted_bigrams = [bigrams[i] for i in sorted_indices]
sorted_values = [values[i] for i in sorted_indices]

# Create a bar chart
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')
ax.bar(range(len(sorted_values)), sorted_values, align='center', color='#00AF87')
ax.set_xticks(range(len(sorted_values)))
ax.set_xticklabels([f'{bigram[0]} {bigram[1]}' for bigram in sorted_bigrams], rotation=45)
ax.set_xlabel('Bigrams', color='white', fontsize=13)
ax.set_ylabel('Values', color='white', fontsize=13)
ax.tick_params(axis='y', colors='white')
ax.set_title('Top 20 Bigrams: Not Detractors', color='white', fontsize=15)
plt.tight_layout()
plt.show()
```



```
In [130]: # Searching through actual reviews to define the themes of some bigrams
# X_train_detractors[X_train_detractors['Review_prep_nolist'].str.contains('cred')]
```

Despite detractors giving overall negative reviews, they also highlight what they liked in the hotel as a contrast. This is why location is often referred to as something great.

I will remove the references to location, as well as the indications such as 'make sure' so I can clearly display the top 10 themes to focus on.

```
In [131]: bigrams_to_remove = [('great', 'location'), ('make', 'sure'), ('good', 'location')]
# Plotting the detractor bigrams as a bar chart

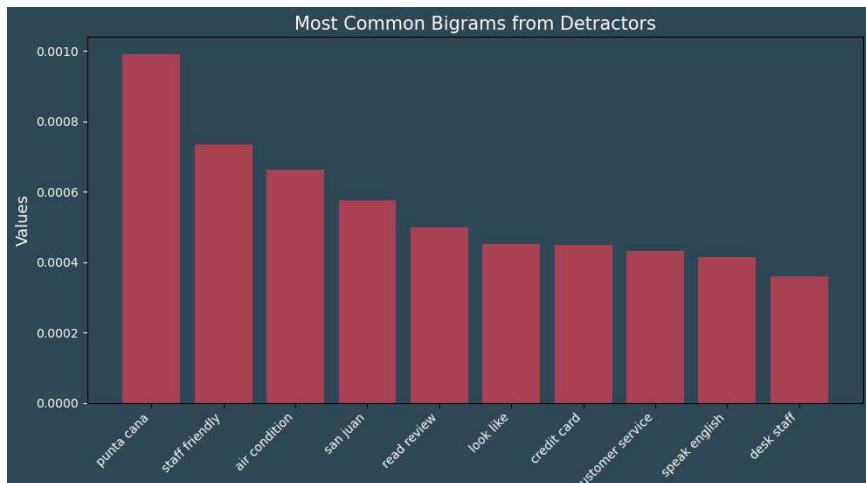
In [132]: # Extracting bigrams and values they are not in list of bigrams to remove
bigrams_for_bar = bigrams_10(*detractors_bigrams_10) detractors_bigrams_500 if bigram[0] in bigrams_to_remove else None
detractors_bigrams_10

# Sort bigrams and values in descending order based on values
sorted_indices = sorted(range(0, len(detractors_bigrams_10)), key=lambda k: values[k], reverse=True)
sorted_bigrams = [bigrams_for_bar[i] for i in sorted_indices]
sorted_values = [values[i] for i in sorted_indices]

# Create a bar chart
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')
ax.bar(range(0, len(sorted_values)), sorted_values, align='center', color='red')
ax.set_xticks(range(0, len(sorted_values)))
ax.set_xticklabels([f'{bigram[0]} {bigram[1]}' for bigram in sorted_bigrams], rotation=45)
ax.set_xlabel('Bigrams', color='white', fontsize=13)
ax.set_ylabel('Values', color='white', fontsize=13)
ax.tick_params(axis='y', colors='white')
ax.set_title('Most Common Bigrams from Detractors', color='white', fontsize=15)
```

```
In [131]: # Plotting the detractor bigrams as a bar chart
```

```
In [132]: # Extracting bigrams and values key are not in List of bigrams to remove
detractors_bigrams_20(*detractors_bigrams_20)
detractors_bigrams_10 = [bigram for bigram in detractors_bigrams_20 if bigram[0] in detractors_bigrams_10]
# Sort bigrams and values in descending order based on values
sorted_indices = sorted(range(len(values)), key=lambda k: values[k], reverse=True)
sorted_bigrams = [bigrams[i] for i in sorted_indices]
sorted_values = [values[i] for i in sorted_indices]
(('san', 'juan'), 0.0005768260148532699),
((('read', 'a book'), 0.0004975124378109452),
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')
ax.bar(range(len(sorted_values)), sorted_values, align='center', color="#AA4255")
ax.set_xticklabels([f'{bigram[0]} {bigram[1]}' for bigram in sorted_bigrams], rotation=45)
ax.set_xlabel('Bigrams', color='white', fontsize=13)
ax.set_ylabel('Values', color='white', fontsize=13)
ax.tick_params(axis='y', colors='white')
ax.set_title('Most Common Bigrams from Detractors', color='white', fontsize=15)
plt.tight_layout()
plt.savefig('images/top_10_bigrams.png')
plt.show()
```



```
In [134]: # Plotting the not_detractor bigrams as a bar chart
```

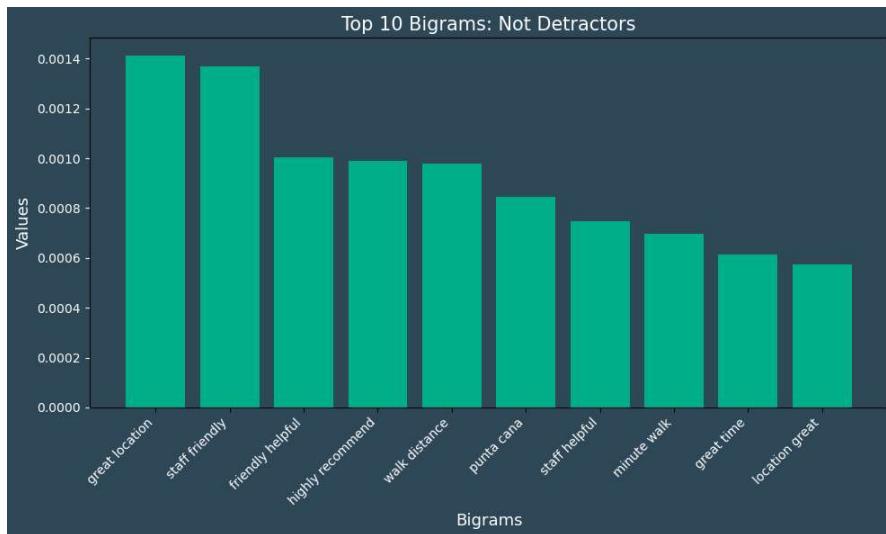
```
# Extracting bigrams and values
bigrams, values = zip(*notdetractors_bigrams_20[:10])
# Sort bigrams and values in descending order based on values
sorted_indices = sorted(range(len(values)), key=lambda k: values[k], reverse=True)
sorted_bigrams = [bigrams[i] for i in sorted_indices]
sorted_values = [values[i] for i in sorted_indices]
# Create a bar chart
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')
ax.bar(range(len(sorted_values)), sorted_values, align='center', color="#00AF87")
ax.set_xticklabels([f'{bigram[0]} {bigram[1]}' for bigram in sorted_bigrams], rotation=45)
ax.set_xlabel('Bigrams', color='white', fontsize=13)
ax.set_ylabel('Values', color='white', fontsize=13)
ax.tick_params(axis='y', colors='white')
ax.set_title('Top 10 Bigrams: Not Detractors', color='white', fontsize=15)
```

```
In [134]: # Plotting the not_detractor bigrams as a bar chart

# Extracting bigrams and values
bigrams, values = zip(*notdetractors_bigrams_20[:10])

# Sort bigrams and values in descending order based on values
sorted_indices = sorted(range(len(values)), key=lambda k: values[k], reverse=True)
sorted_bigrams = [bigrams[i] for i in sorted_indices]
sorted_values = [values[i] for i in sorted_indices]

# Create a bar chart
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#F0F0F0')
ax.set_facecolor('#F0F0F0')
ax.bar(range(len(sorted_values)), sorted_values, align='center', color='#00AF87')
ax.set_xticks(range(len(sorted_values)))
ax.set_xticklabels(['f'{bigram[0]} {bigram[1]}' for bigram in sorted_bigrams], rotation=90)
ax.set_xlabel('Bigrams', color='white', fontsize=13)
ax.set_ylabel('Values', color='white', fontsize=13)
ax.tick_params(axis='y', colors='white')
ax.set_title('Top 10 Bigrams: Not Detractors', color='white', fontsize=15)
plt.tight_layout()
plt.show()
```



Bigrams from detractors indicate that:

- Punta Cana and San Juan are the first destinations that require attention
  - Staff friendliness needs to be addressed as identified in wordclouds
  - Air conditioning needs to work and remain quiet
  - Appearance matters: ensure the hotel is well-maintained
  - Hire staff who can speak English in destinations welcoming English-speaking guests
  - Ensure easy payment methods are allowed and work correctly

Some hotels (i.e. Florence) do not take credit cards. This needs to be addressed nowadays.

In [136]: • Respond to reviews: a lot of detractors mention they read them  
# Calling the function for the top 10 mutual information scores for all reviews

```
Out[136]: 4. (2,1) Mutual Information Scores
[('ajili', 'mojili'), 18.094979706748987),
 ('krispy', 'kreml'), 18.094979706748987),
 ('desa', 'seni'), 17.83194530091519),
 ('esoh', 'pins'), 17.83194530091519)
I will calculate mutual information scores and I will create a frequency filter, so that I only examine
the strength of association between pairs of words (bigrams) for detractors or not detractors
[('ropa', 'vieja'), 17.609552879578743),
 ('altos', 'chavon'), 17.41690780163635),
 ('das', 'ichi'), 17.41690780163635),
 ('kwai', 'fong'), 17.41690780163635),
```

```
In [135]: #(Upgrading a function to edit tf-idf scores)
def mutual_info_score(text_rmi_index = BigramCollocationFinder.from_words(text.sum()))
```

```
In [137]: # Rehashing the pmi to apply frequency filters by detractors or not detractors review
text_pmi_finder = BigramCollocationFinder.from_words(text.sum())
text_pmi_finder.apply_freq_filter(0.05)
X_train_detractors = X_train[X_train['label'] == 'negatives'].apply(bigram_measures.pmi)
X_train_no_detractors = X_train[X_train['label'] != 'detractors'].apply(bigram_measures.pmi)
```

```
In [138]: # Calling the function for the top 10 mutual information scores for positive reviews
mutual_info_score(X_train_detractors['Review'], 5)[:10]
```

In [136]: • Respond to reviews: a lot of detractors mention they read them  
*# Calling the function for the top 10 mutual information scores for all reviews*  
`mutual_info_score(X_train['Review'], 5)[:10]`

Out[136]: **4.2.1 Mutual Information Scores**  
`[('aiili', 'majili'), 18.094979706748987),  
 ('krispy', 'kreme'), 18.094979706748987),  
 ('desa', 'seni'), 17.83194530091519),  
 ('will', 'escolar'), 17.83194530091519),  
 I will calculate mutual information scores and it will create a frequency filter, so that I only examine  
 the strength of association between pairs of words (bigrams) for detractors or not detractors  
 ('altos', 'chavon'), 17.41690780163635),  
 ('sentidos', 'dal'), 17.41690780163635),  
 ('(dal', 'ichi'), 17.41690780163635),  
 ('kwai', 'fong'), 17.41690780163635),`

In [135]: *# Assigning a review to 17.41690780163635 mutual information scores*  
`def mutual_info_score(text):`  
 `text_pmi_finder = BigramCollocationFinder.from_words(text.sum())`

In [137]: *# Removing the previously created subset by detractors or not detractors review*  
`X_train_detractors = X_train[X_train['label'] == 'detractors']`  
`X_train_notdetractors = X_train[X_train['label'] != 'detractors']`

In [138]: *# Calling the function for the top 10 mutual information scores for positive reviews*  
`mutual_info_score(X_train_detractors['Review'], 5)[:10]`

Out[138]: `[('sha', 'tsui'), 16.344538648678256),  
 ('tsim', 'sha'), 16.344538648678256),  
 ('degli', 'orafi'), 16.08150424284446),  
 ('lo', 'behold'), 16.08150424284446),  
 ('rick', 'steves'), 15.859111821508012),  
 ('ograde', 'oreview'), 15.666466743565618),  
 ('chateau', 'lemonye'), 15.496541742123306),  
 ('postage', 'stamp'), 15.496541742123306),  
 ('sagrada', 'familia'), 15.47382166562322),  
 ('zona', 'rosa'), 15.444074322229168)]`

In [139]: *# Calling the function for the top 10 mutual information scores for notdetractors*  
`mutual_info_score(X_train_notdetractors['Review'], 5)[:10]`

Out[139]: `[('kwai', 'fong'), 17.58614342342153),  
 ('desa', 'seni'), 17.323109017587733),  
 ('smith', 'wollensky'), 17.323109017587733),  
 ('piet', 'hein'), 17.10071659625129),  
 ('degli', 'orafi'), 16.908071518308894),  
 ('pont', 'neuf'), 16.908071518308894),  
 ('dai', 'ichi'), 16.908071518308894),  
 ('zoologischer', 'garten'), 16.908071518308894),  
 ('crabtree', 'evelyn'), 16.73814651686658),  
 ('frankfurter', 'allee'), 16.685679096972443)]`

The mutual information scores in themselves are a bit harder to drive direct conclusions from them. Additional analyses would be necessary for them to be conclusive.

The groups of words mostly indicate some specific neighborhoods in destinations that were not appreciated by guests, or at the opposite, some destinations and attractions in destinations that were often talked by guests who reviewed the hotels positively.

## • Term-Frequency

**5. Modeling** often a term (word) appears in a document

- Inverse Document Frequency (IDF)

Measures the importance of a term in the entire collection of documents.

My objective is now to:

/ main classification models were explored:

1. Develop a tool to identify unsatisfied customer reviews in real time
2. Multinomial Naive Bayes

Incorrectly identifying a guest as satisfied when they weren't would lead in a detractor posting a review instead of catching them while they are in-house, and having a chance of improving their satisfaction.

3. Decision Tree

As a consequence, recall is the main evaluation metric for this model.

4. Random Forest

5. Gradient Boosting

As the dataset is a text, it requires a transformation before it can be used for modeling. Like other Undersampling, stopwords, lemmatize, and hyperparameter tuning were parameters used to tune types of dataset would one-hot encoded, here, the reviews were vectorized using the common models. In addition, X\_train was split with `stratify` in order to keep the distribution ratios for each method in natural language processing: TfidfVectorizer .

It converts a collection of text documents to a matrix of tf-idf features.

- Term-Frequency

## 5. Modeling

Measures how often a term (word) appears in a document

- Inverse Document Frequency (IDF)

I now have an initial idea for recommendations on focus areas to improve guest satisfaction.

My objective is now to:

7 main classification models were explored:

1. Develop a tool to identify unsatisfied customer reviews in real time

### 1. Multinomial Naive Bayes

Incorrectly labeling a guest as satisfied when they weren't would lead in a detractor posting a review instead of catching them while they are in-house, and having a chance of improving their satisfaction.

### 2. K-Nearest Neighbor

Incorporating a guest as satisfied when they weren't would lead in a detractor posting a review instead of catching them while they are in-house, and having a chance of improving their satisfaction.

### 3. Decision Tree

Incorporating a guest as satisfied when they weren't would lead in a detractor posting a review instead of catching them while they are in-house, and having a chance of improving their satisfaction.

### 4. Random Forest

Incorporating a guest as satisfied when they weren't would lead in a detractor posting a review instead of catching them while they are in-house, and having a chance of improving their satisfaction.

### 5. Gradient Boosting

As a consequence, recall is the main evaluation metric for this model.

### 6. XGBoost

As a consequence, recall is the main evaluation metric for this model.

As a consequence, recall is the main evaluation metric for this model.

### 7. XGBoost

As the dataset is a text, it requires a transformation before it can be used for modeling. Like other

Undersampling, stopwords, lemmatize, and hyperparameter tuning were parameters used to tune

types of dataset would one-hot encoded. here, the reviews were vectorized, using the common

models. In addition, `X_train` was split with `stratify` in order to keep the distribution ratios for each

method in natural language processing: `TfidfVectorizer`.

sentiment.

It converts a collection of text documents to a matrix of tf-idf features.

```
In [140]: # Reminding the natural balance
y_train.value_counts(normalize=True)
```

```
Out[140]: 0    0.736596
1    0.263404
Name: Sentiment, dtype: float64
```

If I were to guess the majority class every time I would get 74% accuracy.

```
In [141]: # Verifying Labels is not part of X_train
X_train.columns
```

```
Out[141]: Index(['Review', 'Rating', 'label', 'original_review', 'Review_prep_nolist'], d
type='object')
```

```
In [142]: # Dropping it
# X_train = X_train.drop('Label', axis=1)
```

By using only the dataset's natural class balance, and if I guessed the contribution of the majority class every time I would get 74% accuracy. However, if I were to guess that a review was not from a detractor, I would expect only about 26% accuracy.

## 5. a) Baseline Model with TfidfVectorizer and MultinomialNB

The first baseline model will vectorize the reviews and make predictions using Multinomial Naive Bayes. The first step is to import the vectorizer, instantiate a vectorizer object and fit it on `X_train['original_review']`.

### 1st iteration: Vanilla Tfidf Vectorizer with Pipeline

```
In [143]: # Importing the relevant packages
from sklearn.pipeline import Pipeline
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.naive_bayes import MultinomialNB

# Define the pipeline steps
tfidf_vectorizer = TfidfVectorizer()
naive_bayes_classifier = MultinomialNB()

# Create the pipeline
base_pipeline = Pipeline([
    ('tfidf', tfidf_vectorizer),
    ('classifier', naive_bayes_classifier)
])

# Fitting the pipeline on X_train['original_review'] and y_train
base_pipeline.fit(X_train['original_review'], y_train)

# Calculating predictions using this model
base_y_pred = base_pipeline.predict(X_test['original_review'])
```

```
In [143]: # Importing the relevant packages
from sklearn.pipeline import Pipeline
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.naive_bayes import MultinomialNB

# Define the pipeline steps
tfidf_vectorizer = TfidfVectorizer()
naive_bayes_classifier = MultinomialNB()

# Create the pipeline
base_pipeline = Pipeline([
    ('tfidf', tfidf_vectorizer),
    ('classifier', naive_bayes_classifier)
])

# Fitting the pipeline on X_train['original_review'] and y_train
base_pipeline.fit(X_train['original_review'], y_train)

# Calculating predictions using this model
base_y_pred = base_pipeline.predict(X_test['original_review'])

# Optionally, you can access the individual components of the pipeline:
X_train_vectorized = base_pipeline.named_steps['tfidf'].transform(X_train['original_review'])
baseline_model = base_pipeline.named_steps['classifier']
```

## 2) Evaluation Metrics

```
In [144]: # Importing the relevant packages
from sklearn.metrics import recall_score, accuracy_score, f1_score
from sklearn.model_selection import cross_val_score

def evaluation_metrics(y_test, y_pred, model, X, y):
    with warnings.catch_warnings():
        warnings.simplefilter("ignore", category=RuntimeWarning)

    # Calculating and printing accuracy
    accuracy = accuracy_score(y_test, y_pred)
    print(f'Accuracy: {accuracy:.4f}')

    # Calculating and printing F1-score
    f1 = f1_score(y_test, y_pred, average='weighted')
    print(f'F1-Score: {f1:.4f}')

    # Calculating and printing recall
    recall = recall_score(y_test, y_pred, average='weighted', zero_division=0)
    print(f'Recall: {recall:.4f}')

    # Performing cross-validation and printing the mean accuracy
    cv_scores = cross_val_score(model, X, y, cv=5)
    mean_cv_accuracy = cv_scores.mean()
    print(f'Mean Cross-Validated Accuracy: {mean_cv_accuracy:.4f}')

    return accuracy, f1, recall, mean_cv_accuracy
```

```
In [145]: # Naming the model and calling the function to evaluate it
baseline_model_name = 'Baseline'
    3) Classification Report
# Calling the function and recording into the defined values
accuracy_base, f1_base, recall_base, cv_base = evaluation_metrics(
    y_test,
    base_y_pred,
    base_pipeline,
    import_X_train['original_review'],
    from sklearn.metrics import classification_report
    y_train)

def plays_oakley(y_test, y_pred):
    F1-Score: 0.6656
    Recall: 0.7546
    Mean Cross-Validated Accuracy: 0.7500
    # Cross-Validated Accuracy: classification report
    with warnings.catch_warnings():
        warnings.simplefilter("ignore", category=RuntimeWarning)
The model shows a reasonable overall accuracy in predicting sentiments from hotel reviews on TripAdvisor. It achieves a good balance between recall and F1-Score, with a mean cross-validated accuracy that suggests stability across different subsets of the data.
    # Due to class imbalance, the initial recall for positive returns 0
    # Consequently, warnings need to be handled
```

```
In [145]: # Naming the model and calling the function to evaluate it
baseline_model_name = 'Baseline'
    3) Classification Report
# Calling the function and recording into the defined values
accuracy_base, f1_base, recall_base, cv_base = evaluation_metrics(
In [146]: # Defining a function to print a classification report
base_y_pred,
base_pipeline,
import warnings
X_train['Original_review'],
from sklearn.metrics import classification_report
y_train)

def class_calculation(y_test, y_pred):
    # F1-Score will be calculated for each model beforehand
    Recall: 0.7546
    Mean # Cross-Validated Accuracy: 0.7502
    with warnings.catch_warnings():
        warnings.simplefilter("ignore", category=RuntimeWarning)
The model shows a reasonable overall accuracy in predicting sentiments from hotel reviews on TripAdvisor achieves a good balance between recall and F1-Score with the mean cross-validated
accuracy that suggests stability across different subsets of the data.
    # Due to class imbalance, the initial recall for positive returns 0
    # Consequently, warnings need to be handled
    print('Classification Report:\n', class_report)

    # Parsing the classification report to extract recall for 'Detractors'
lines = class_report.split('\n')
recall_line = lines[3]
# Extracting numerical values
recall_values = [float(value) for value in recall_line.split()[1:]]
# Storing recall for Detractors
recall_detractors = recall_values[1]

return class_report, recall_detractors
```



Because what matters most is recall but for the detractors' side, we will also store this metric.

```
In [147]: # Calling the function to record the classification report
base_class_report, recall_detractors_base = class_calculation(y_test, base_y_pred)

Classification Report:
precision      recall   f1-score   support
Not Detractors  0.75020  0.99973  0.85718   3773
    Detractors   0.98947  0.06963  0.13010   1350

    accuracy       0.75464
    macro avg     0.86984  0.53468  0.49364   5123
weighted avg    0.81325  0.75464  0.66558   5123
```

```
In [149]: # Display confusion matrix the confusion matrix

# Importing the relevant package
from sklearn.metrics import confusion_matrix
from sklearn.metrics import ConfusionMatrixDisplay
from matplotlib.colors import LinearSegmentedColormap

def confusion_matrix_display(model, y_test, y_pred):
    # Defining a colormap that interpolates between the defined colors
    custom_colors_cnf = ['#00AF87', '#00A4A0', '#0097B9', '#0097B9']
    cnf_matrix = confusion_matrix(y_test, y_pred)
    # print(cnf_matrix)
    custom_colors_detractors = ['#F3E8EA', '#E5BEC4', '#AA4255']
    # custom_colors_detractors_reversed = ['#AA4255', '#E5BEC4', '#F3E8EA']
    # Normalizing the confusion matrix
    cnf_matrix_normalized = cnf_matrix.astype('float') / cnf_matrix.sum(axis=1)[
    n_bins = 3
    # Creating in the model's classes as labels
    disp = ConfusionMatrixDisplay(confusion_matrix=cnf_matrix_normalized, display_labels=custom_cmap_cnf)
    # Taking in pre-defined labels
    disp = ConfusionMatrixDisplay(confusion_matrix=cnf_matrix_normalized, display_labels=custom_cmap_cnf)
    disp.plot(cmap=custom_cmap_cnf)
```

```
In [149]: # DisplayConfusionMatrix the confusion matrix

# Importing the relevant package
from sklearn.metrics import confusion_matrix
from sklearn.metrics import ConfusionMatrixDisplay
from matplotlib.colors import LinearSegmentedColorMap

def confusion_matrix_display(model, y_test, y_pred):
    # Defining a colormap that interpolates between the defined colors
    custom_colors_cnf = ['#00AFA0', '#004AA0', '#0097B9', '#0097B9']
    cnf_matrix = confusion_matrix(y_test, y_pred)
    print(cnf_matrix)
    custom_colors_detractors = ['#F3E8EA', '#E5BEC4', '#AA4255']
    # custom_colors_detractors_reversed = ['#AA4255', '#E5BEC4', '#F3E8EA']
    # Normalizing the confusion matrix
    cnf_matrix_normalized = cnf_matrix.astype('float') / cnf_matrix.sum(axis=1)[
        :, None]
    n_bins = 5
    # Taking in the model's classes as labels
    # disp = ConfusionMatrixDisplay(confusion_matrix=cnf_matrix_normalized, display_
    #                                 values=True)
    custom_cmap_cnf = LinearSegmentedColorMap.from_list("custom_cmap_cnf", custom_color
    # Taking in pre-defined labels
    disp = ConfusionMatrixDisplay(confusion_matrix=cnf_matrix_normalized, display_
        values=False)
    disp.plot(cmap=custom_cmap_cnf)

    plt.title("Model Performance: Confusion Matrix", fontsize=16)

    # Saving the plot as a PNG with a transparent background
    # plt.savefig('images/confusion_matrix.png', transparent=True)

    return cnf_matrix
    plt.show()
```

```
In [150]: # Creating a copy for the best model, with white values for labels and titles

def best_confusion_matrix_only(model, y_test, y_pred):
    # Defining the confusion matrix
    cnf_matrix = confusion_matrix(y_test, y_pred)

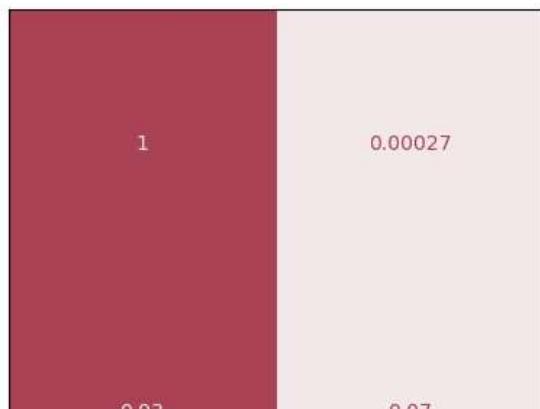
    # Normalizing the confusion matrix
    cnf_matrix_normalized = cnf_matrix.astype('float') / cnf_matrix.sum(axis=1)[
        :, None]
    print(cnf_matrix_normalized)

    # Taking in pre-defined labels
    disp = ConfusionMatrixDisplay(confusion_matrix=cnf_matrix_normalized, display_
        values=False)
    disp.plot(cmap=custom_cmap_cnf, colorbar=False) # Set colorbar to False to remove it

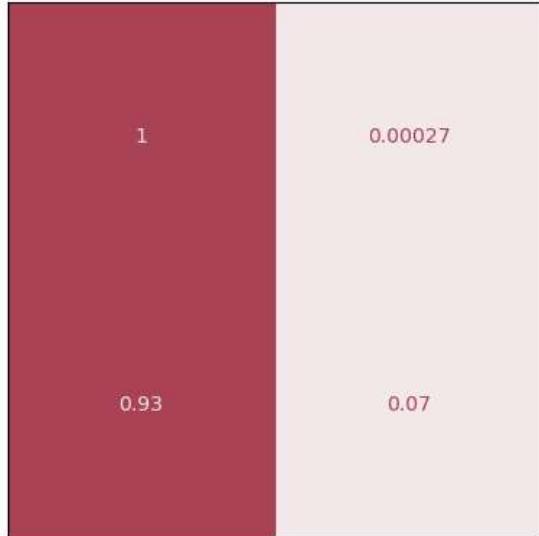
    # Customizing the plot
    # Setting the title and labels color to white
    plt.title("Model Performance: Confusion Matrix", fontsize=16, color='white')
    # Setting the color of x-axis and y-axis labels to white

    disp.ax_.xaxis.label.set_color("white")
    disp.ax_.yaxis.label.set_color("white")
    disp.ax_.tick_params(axis='x', colors='white')
    disp.ax_.tick_params(axis='y', colors='white')

    # Saving the plot as a PNG with a transparent background
    # plt.savefig('images/best_confusion_matrix.png', transparent=True)
    best_confusion_matrix_only(base_pipeline, y_test, base_y_pred)
    plt.show()
```

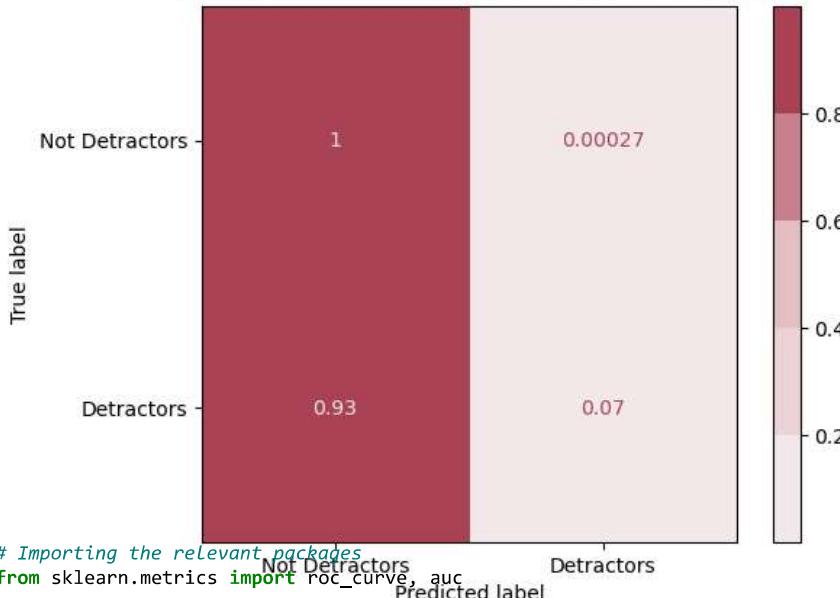


```
# Saving the plot as a PNG with a transparent background
In [151]: # plt.savefig('images/best_confusion_matrix.png', transparent=True)
best_confusion_matrix_only(base_pipeline, y_test, base_y_pred)
plt.show()
```



```
# Calling the function to display the confusion matrix
In [152]: base_cfn_matrix = confusion_matrix_display(base_pipeline, y_test, base_y_pred)
```

Model Performance: Confusion Matrix



```
# Importing the relevant packages
In [153]: from sklearn.metrics import roc_curve, auc

def plot_roc_curve(y_true, y_pred_proba, model_name='Model'):
    # Computing ROC curve and AUC
    fpr, tpr, thresholds = roc_curve(y_true, y_pred_proba)
    roc_auc = auc(fpr, tpr)

    # Plotting ROC curve
    plt.figure(figsize=(8, 6))
    plt.plot(fpr, tpr, color="#AA4255", lw=2, label='{} (AUC = {:.4f})'.format(model_name, roc_auc))

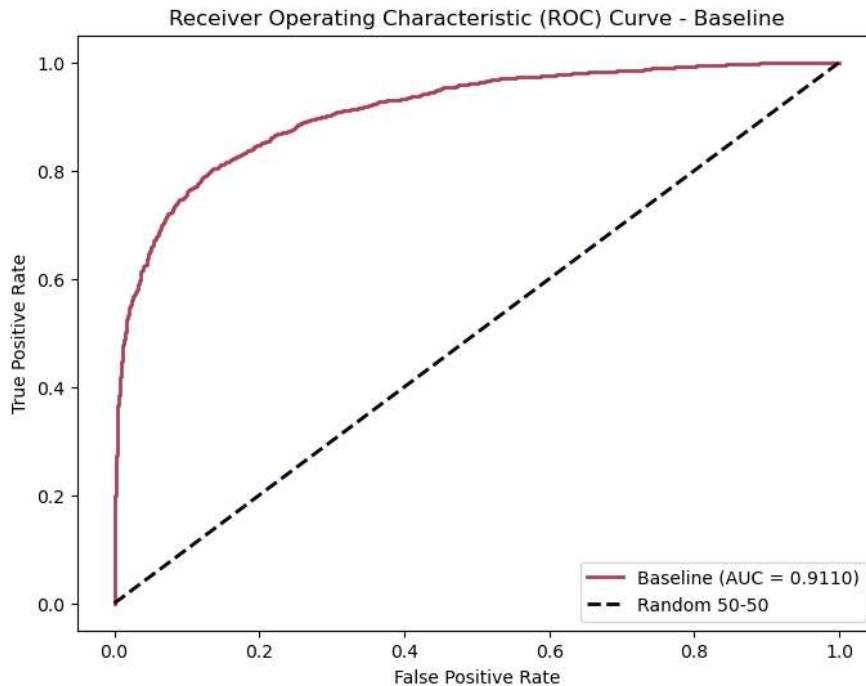
    # Plotting the 50-50 accuracy line
    plt.plot([0, 1], [0, 1], linestyle='--', color='black', label='Random')
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.title('Receiver Operating Characteristic (ROC) Curve - {}'.format(model_name))
    plt.legend(loc='lower right')

    plt.show()
    return roc_auc
```

The Receiver Operating Characteristic (ROC) curve displays a graph where the True Positive Rate against the False Positive Rate. It is another evaluation metric for classification algorithms combining both described ratios. It is compared to a *worthless accuracy*: a classifier with 50-50 accuracy. It is considered *worthless*, as it would be comparable to just random guessing.

```
In [153]: # Importing the relevant packages
from sklearn.metrics import roc_curve, auc
# Computing ROC curve and AUC
def plot_roc_curve(y_true, y_pred_proba, model_name='Model'):
    While the overall metrics took decent, the results for each sentiment are not balanced at all and the
    recall for detractors is the lowest of recorded results (0.0696). This needs to be addressed through
    the dataset balance of fpr, tpr)
    # Plotting ROC curve
    plt.figure(figsize=(8, 6))
    plt.plot(fpr, tpr, color="#AA4255", lw=2, label='{} (AUC = {:.4f})'.format(model_name, auc))
    # Plotting the 50-50 accuracy line
    The Receiver Operating Characteristic (ROC) Curve displays 2 lines: a solid line for the Model's Performance
    against the True Positive Rate. It is another evaluation metric for classification algorithms
    combining both described ratios. It is compared to a worthless accuracy: a classifier with 50-50
    accuracy is considered worthless as it would be comparable to just random guessing.
    plt.title('Receiver Operating Characteristic (ROC) Curve - {}'.format(model_name))
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.legend(loc='lower right')
    plt.show()
    return roc_auc
```

```
In [154]: # Plotting the ROC curve for the baseline model
base_y_pred_proba = base_pipeline.predict_proba(X_test['original_review'])[:, 1]
auc_base = plot_roc_curve(y_test, base_y_pred_proba, model_name=baseline_model_name)
```



ROC-AUC is a common classification metric for binary classification problems to evaluate the ability of a model to distinguish between the positive and negative classes. It visually represents the trade-off between the True Positive rate and the False Positive rate at various thresholds. The area under the ROC curve (AUC) provides a single value that summarizes the model's performance. The higher the AUC is, the better the model can distinguish between the classes.

```
# Defining a function that provides an interpretation of area under the curve
def auc_interpretation(auc):
    if auc == 1.0:
        print(f'AUC = {auc:.4f}. The model is a perfect classifier.')
    elif auc >= 0.95:
        print(f'AUC = {auc:.4f}. The model has a predictive power close to perfect.')
    elif auc > 0.9:
        print(f'AUC = {auc:.4f}. The model has a great predictive power in distinguishing between classes.')
    elif auc > 0.8:
        print(f'AUC = {auc:.4f}. The model has a good predictive power.')
    elif auc > 0.7:
        print(f'AUC = {auc:.4f}. The model has some predictive power.')
    elif auc == 0.5:
        print(f'AUC = {auc:.4f}. The model does not perform better than random chance.')
```

```
In [156]: # Printing the interpretation
auc_interpretation(auc_base)
```

In [155]: # Defining the AUC function. This prints the interpretation of the area under the ROC curve.

The area under the ROC curve (AUC) provides a single value that summarizes the model's performance. The higher the AUC is, the better the model can distinguish between the classes.

```

def print_auc(auc):
    if auc == 1.0:
        print(f'AUC = {auc:.4f}. The model is a perfect classifier.')
    elif auc >= 0.95:
        print(f'AUC = {auc:.4f}. The model has a predictive power close to perfect.')
    elif auc > 0.9:
        print(f'AUC = {auc:.4f}. The model has a great predictive power in distinguishing between the positive and negative classes.')
    elif auc > 0.8:
        print(f'AUC = {auc:.4f}. The model has a good predictive power.')
    elif auc > 0.7:
        print(f'AUC = {auc:.4f}. The model has some predictive power.')
    elif auc == 0.5:
        print(f'AUC = {auc:.4f}. The model does not perform better than random chance.')
    else:
        print(f'AUC = {auc:.4f}. The model has a low predictive power in distinguishing between the positive and negative classes.')

```

In [156]: # Printing the interpretation

```
auc_interpretation(auc_base)
```

AUC = 0.9110. The model has a great predictive power in distinguishing between the positive and negative classes.

In [157]: # Creating an empty list of metrics to easily add new evaluations as they are done.

```

modelnames = []
accuracies = []
f1s = []
recalls = []
recalls_detractors = []
cv_acccs = []
roc_aucs = []

```

In [158]: # Appending each metric to the lists

```

modelnames.append(baseline_model_name)
accuracies.append(accuracy_base)
f1s.append(f1_base)
recalls.append(recall_base)
recalls_detractors.append(recall_detractors_base)
cv_acccs.append(cv_base)
roc_aucs.append(auc_base)

```

## 2nd iteration: Addressing class imbalance: undersampling negative reviews

### Sampling Strategy: 0.8

While the initial categorization of sentiment following Net Promoter Score classification helped address the dataset imbalance, the number of reviews from non-detectors remained higher than the one from detractors.

As a consequence, the dataset needs to be resampled. More precisely, negative reviews need to

In [160]: # Defining the pipeline steps

```

tfidf_vectorizer = TfidfVectorizer()
naive_bayes_classifier = MultinomialNB()
    1) Fitting and training on train data
# Including the UnderSampler to the pipeline
rs_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', tfidf_vectorizer),
    ('imblearn.RandomUnderSampler(sampling_strategy=0.8, random_state=42)),
    ('classifier', naive_bayes_classifier)
])

# will now include the RandomUnderSampler to our pipeline so it undersamples our majority class.
rs_pipeline.fit(X_train['original_review'], y_train)

rs_y_pred = rs_pipeline.predict(X_test['original_review'])

```

In [161]: # Verifying the value counts of each category, before undersampling

```

original_value_counts = y_train.value_counts()
original_value_counts_normalized = y_train.value_counts(normalize=True)
print("Original class distribution:")

```

```
In [160]: be undersampled. I will first try with a sampling strategy of 0.8.
# Defining the pipeline steps
tfidf_vectorizer = TfidfVectorizer()
naive_bayes_classifier = MultinomialNB()
    1) Fitting and training on train data
# Including the UnderSampler to the pipeline
rs_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', tfidf_vectorizer),
    # Importing relevant packages
    from imblearn.under_sampling import RandomUnderSampler(sampling_strategy=0.8, random_state=42),
    ('classifier', naive_bayes_classifier)
])

In [159]: # Will now include the Random UnderSampler to our pipeline so it undersamples our majority class.
rs_pipeline.fit(X_train['original_review'], y_train)

rs_y_pred = rs_pipeline.predict(X_test['original_review'])
```

```
In [161]: # Verifying the value counts of each category, before undersampling
original_value_counts = y_train.value_counts()
original_value_counts_normalized = y_train.value_counts(normalize=True)
print("Original class distribution:")
print(original_value_counts)
print(original_value_counts_normalized)

# Getting the indices of the resampled data
resampled_indices = rs_pipeline.named_steps['us'].sample_indices_

# Verifying the new value counts of each category, after undersampling
resampled_value_counts = y_train.iloc[resampled_indices].value_counts()
resampled_value_counts_normalized = y_train.iloc[resampled_indices].value_counts(normalize=True)
print("\nClass distribution after undersampling:")
print(resampled_value_counts)
print(resampled_value_counts_normalized)
```

```
Original class distribution:
0    11320
1     4048
Name: Sentiment, dtype: int64
0    0.736596
1    0.263404
Name: Sentiment, dtype: float64

Class distribution after undersampling:
0     5060
1     4048
Name: Sentiment, dtype: int64
0    0.555556
1    0.444444
Name: Sentiment, dtype: float64
```

The non-detectors reviews were randomly undersampled from our model, which now has a higher proportion of detractors reviews and lower not\_detractor ones.

I will now run our evaluation metrics to review if the model is providing better predictions.

```
In [163]: # Creating a function to print the metrics' comparison
def metric_comparison(metric_name, old_metric, new_metric):
    if new_metric > old_metric:
        print(f'{metric_name} improved from {old_metric:.4f} to {new_metric:.4f}')
    else:
        print(f'{metric_name} decreased from {old_metric:.4f} to {new_metric:.4f}')

In [162]: # Naming the model
resampled_model_name = 'Resampled'

In [164]: # Comparing the evaluation metrics between the base and undersampled model
metrics_base = {'Accuracy', 'Recall', 'Precision', 'F1', 'Cross-validated Accuracy'}
metrics_rs = {'Accuracy', 'Recall', 'Precision', 'F1', 'Cross-validated Accuracy'}
for metric in metrics_base:
    metric_comparison(metric, f1_base, f1_rs)
for metric in metrics_rs:
    metric_comparison(metric, recall_base, recall_rs)
for metric in metrics_rs:
    metric_comparison(metric, cv_base, cv_rs)
Accuity_improved from 0.7546 to 0.8842.
F1 improved from 0.6656 to 0.8784.
Recall improved from 0.7546 to 0.8842.
Cross-validated Accuracy improved from 0.7507 to 0.8836.
Recall: 0.8842
Mean_Cross-Validated_Accuracy: 0.8836

In [165]: # Printing the overall summary
print('The resampling started to improve the overall scores')
```

```
In [163]: # Creating a function to print the metrics' comparison
def metric_evaluation(metric_name, old_metric, new_metric):
    if new_metric > old_metric:
        print(f'{metric_name} improved from {old_metric:.4f} to {new_metric:.4f}')
    else:
        print(f'{metric_name} decreased from {old_metric:.4f} to {new_metric:.4f}')

In [162]: # Naming the model
resampled_model_name = 'Resampled'

In [164]: # Comparing the baseline and resampling methods defined in the resampled model
# Metrics comparison ('Recall', 'Cross-validated Accuracy')
metric_comparison('Recall', recall_base, recall_rs)
metric_comparison('F1', f1_base, f1_rs)
metric_comparison('Recall', recall_base, recall_rs)
metric_comparison('Cross-validated Accuracy', cv_base, cv_rs)

    X_train['original_review'],
Accuracy improved from 0.7546 to 0.8842.
F1 improved from 0.6656 to 0.8784.
Recall improved from 0.7546 to 0.8842.
Cross-validated Accuracy improved from 0.7507 to 0.8836.
Recall: 0.8842
Mean Cross-Validated Accuracy: 0.8836

In [165]: # Printing the overall summary
print('The resampling started to improve the overall scores')
```

The resampling started to improve the overall scores

### 3) Classification Report

```
In [166]: # Calling confusion report for resampled model
resamp_class_report, recall_detractors_rs = class_calculation(y_test, rs_y_pred)

Classification Report:
precision      recall      f1-score     support
Not Detractors  0.88667   0.96634   0.92479    3773
    Detractors   0.87438   0.65481   0.74884    1350

    accuracy          0.88425    5123
    macro avg       0.88053   0.81058   0.83681    5123
    weighted avg    0.88343   0.88425   0.87843    5123
```

### 4) Confusion Matrix

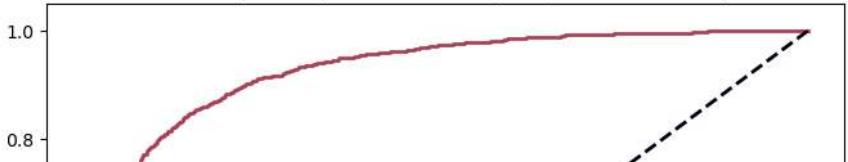
```
In [167]: # Calling confusion matrix for resampled model
rs_cfn_matrix = confusion_matrix_display(rs_pipeline, y_test, rs_y_pred)
```

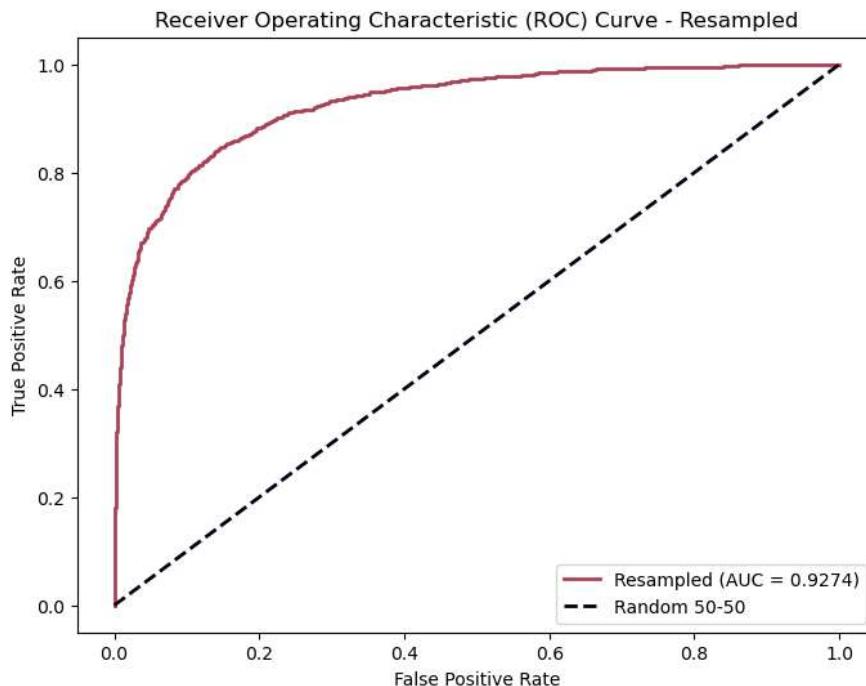
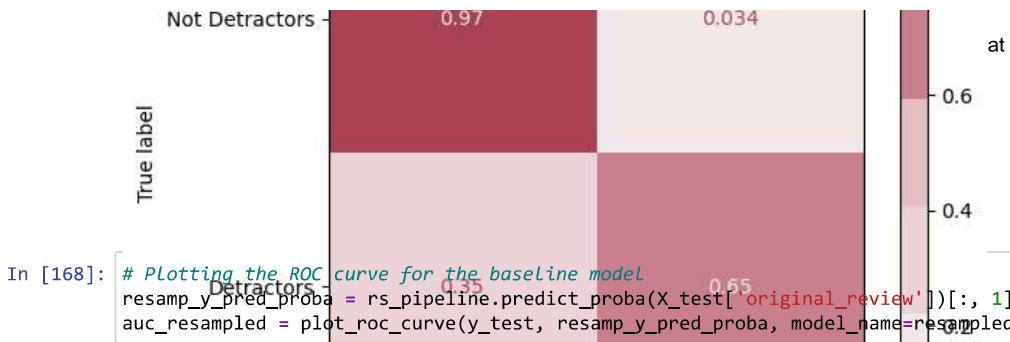
Model Performance: Confusion Matrix



```
In [168]: # Plotting the ROC curve for the baseline model
resamp_y_pred_proba = rs_pipeline.predict_proba(X_test['original_review'])[:, 1]
auc_resampled = plot_roc_curve(y_test, resamp_y_pred_proba, model_name='resampled')
```

Receiver Operating Characteristic (ROC) Curve - Resampled





In [169]:

```
# Printing the interpretation
auc_interpretation(auc_resampled)
```

AUC = 0.9274. The model has a great predictive power in distinguishing between the positive and negative classes.

In [170]:

```
# Comparing AUC to previous one
metrics_comparison('AUC', auc_base, auc_resampled)
```

AUC improved from 0.9110 to 0.9274.

AUC increased, which was to be expected, since the majority class was undersampled to give more weight to the minority class.

In [171]:

```
# Appending each metric to the lists
modelnames.append(resampled_model_name)
accuracies.append(accuracy_rs)
f1s.append(f1_rs)
recalls.append(recall_rs)
```

In [172]:

```
# Importing relevant packages
from imblearn.under_sampling import RandomUnderSampler
recalls_detractors.append(recall_detractors_rs)
cv_acccs.append(cv_ps)
roc_aucs.append(auc_resampled)
```

I will now include the Random Undersampler to our pipeline so it under samples our majority class.

**3rd iteration: Undersampling negative reviews: Sampling Strategy**

In [173]:

```
# Defining the pipeline steps
tfidf_vectorizer = TfidfVectorizer()
naive_bayes_classifier = MultinomialNB()
Sampling Strategy: 1
# Including the UnderSampler to the pipeline
rs1_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', tfidf_vectorizer),
    # Will try to reach the same number of detractors reviews than there are non_detractors reviews.
    # To do so, I will apply a sampling strategy of 1.
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', naive_bayes_classifier)
])
```

```
In [171]: # Appending each metric to the lists
modelnames.append(resampled_model_name)
accuracies.append(accuracy_rs)
f1s.append(f1_rs)

In [172]: # Importing relevant packages
from imblearn.under_sampling import RandomUnderSampler
recall_detractors.append(recall_detractors_rs)
cv_acccs.append(cv_rs)
roc_aucs.append(auc_resampled)
```

I will now include the Random Undersampler to our pipeline so it under samples our majority class.

### 3rd iteration: Undersampling negative reviews: Sampling Strategy

```
In [173]: # Defining the pipeline steps
tfidf_vectorizer = TfidfVectorizer()
naive_bayes_classifier = MultinomialNB()
Sampling Strategy: 1

# Including the UnderSampler to the pipeline
rs1_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', tfidf_vectorizer),
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', naive_bayes_classifier)
])

# Fitting the pipeline on X_train['original_review'] and y_train
rs1_pipeline.fit(X_train['original_review'], y_train)

rs1_y_pred = rs1_pipeline.predict(X_test['original_review'])
```

```
In [174]: # Verifying the value counts of each category, before undersampling
original_value_counts = y_train.value_counts()
print("Original class distribution:")
print(original_value_counts)

# Getting the indices of the resampled data
resampled_indices = rs1_pipeline.named_steps['us'].sample_indices_

# Verifying the new value counts of each category, after undersampling
resampled_value_counts = y_train.iloc[resampled_indices].value_counts()
print("\nClass distribution after undersampling:")
print(resampled_value_counts)
```

Original class distribution:

0	11320
1	4048
Name: Sentiment, dtype: int64	

Class distribution after undersampling:

0	4048
1	4048
Name: Sentiment, dtype: int64	

The non positive reviews were randomly undersampled from our model, which now has the exact same number of reviews for both sentiments.

I will now run our evaluation metrics to review if the model is providing better predictions.

```
In [175]: # Naming the model
resampled1_model_name = 'Sampling 1'

# Evaluation Metrics
# Calling the function and recording into the defined values
accuracy_rs1, f1_rs1, recall_rs1, cv_rs1 = evaluation_metrics(
    y_test,
    rs1_y_pred,
    rs1_pipeline,
    X_train['original_review'],
    y_train)
```

Accuracy: 0.8312  
F1-Score: 0.8378  
Recall: 0.8312  
Mean Cross-Validated Accuracy: 0.8442

```
In [176]: # Comparing the evaluation metrics between baseline and Resampled model
metrics_comparison('Accuracy', accuracy_rs, accuracy_rs1)
metrics_comparison('F1', f1_rs, f1_rs1)
metrics_comparison('Recall', recall_rs, recall_rs1)
metrics_comparison('Cross-validated Accuracy', cv_rs, cv_rs1)
```

I will now run our evaluation metrics to review if the model is providing better predictions.

```
In [175]: # Naming the model
resampled1_model_name = 'Sampling 1'

# Evaluation Metrics
accuracy_rs1, f1_rs1, recall_rs1, cv_rs1 = evaluation_metrics(
    y_test,
    rs1_y_pred,
    rs1_pipeline,
    X_train['original_review'],
    y_train)
```

Accuracy: 0.8312  
F1-Score: 0.8378  
Recall: 0.8312  
Mean Cross-Validated Accuracy: 0.8442

```
In [176]: # Comparing the evaluation metrics between baseline and Resampled model
metrics_comparison('Accuracy', accuracy_rs, accuracy_rs1)
metrics_comparison('F1', f1_rs, f1_rs1)
metrics_comparison('Recall', recall_rs, recall_rs1)
metrics_comparison('Cross-validated Accuracy', cv_rs, cv_rs1)
```

Accuracy decreased from 0.8842 to 0.8312.  
F1 decreased from 0.8784 to 0.8378.  
Recall decreased from 0.8842 to 0.8312.  
Cross-validated Accuracy decreased from 0.8836 to 0.8442.

```
In [177]: # Printing the overall summary
print('The resampling with sampling strategy of 1 did not improve the overall scores')
```

The resampling with sampling strategy of 1 did not improve the overall scores

### 3) Classification Report

```
In [178]: # Calling confusion report for resampled model
resamp1_class_report, recall_detractors_rs1 = class_calculation(y_test, rs1_y_pred)

Classification Report:
precision    recall   f1-score   support
Not Detractors  0.93637  0.82693  0.87825   3773
    Detractors   0.63540  0.84296  0.72461   1350

accuracy          0.83115
macro avg       0.78589  0.83495  0.80143   5123
weighted avg    0.85706  0.83115  0.83777   5123
```

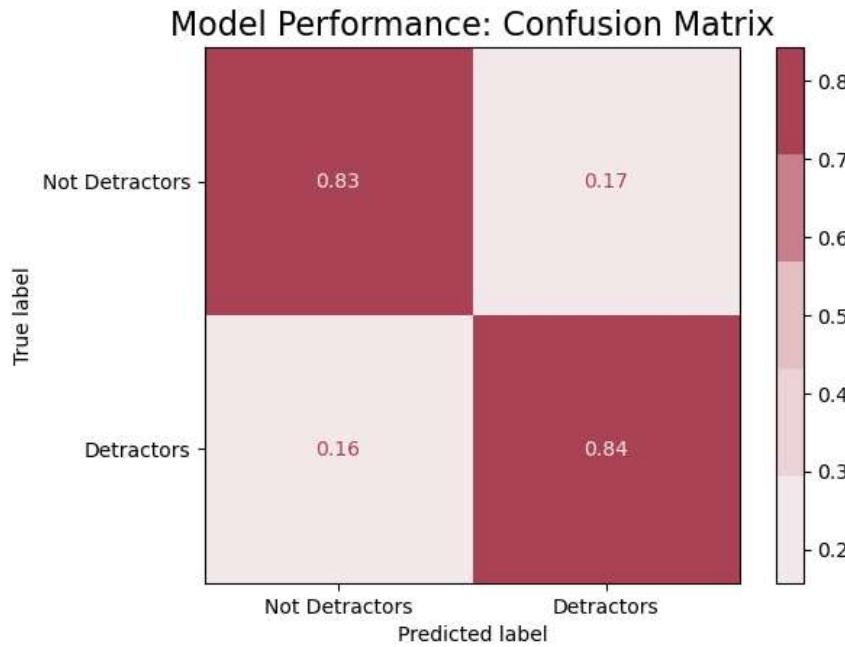
While the overall scores did not improve, the balance between each class: detractors and not detractors, look much better.

Recall for the Detractors class is the highest it was so far.

```
In [179]: # Calling confusion matrix for resampled model
rs1_cfn_matrix = confusion_matrix_display(rs1_pipeline, y_test, rs1_y_pred)
```

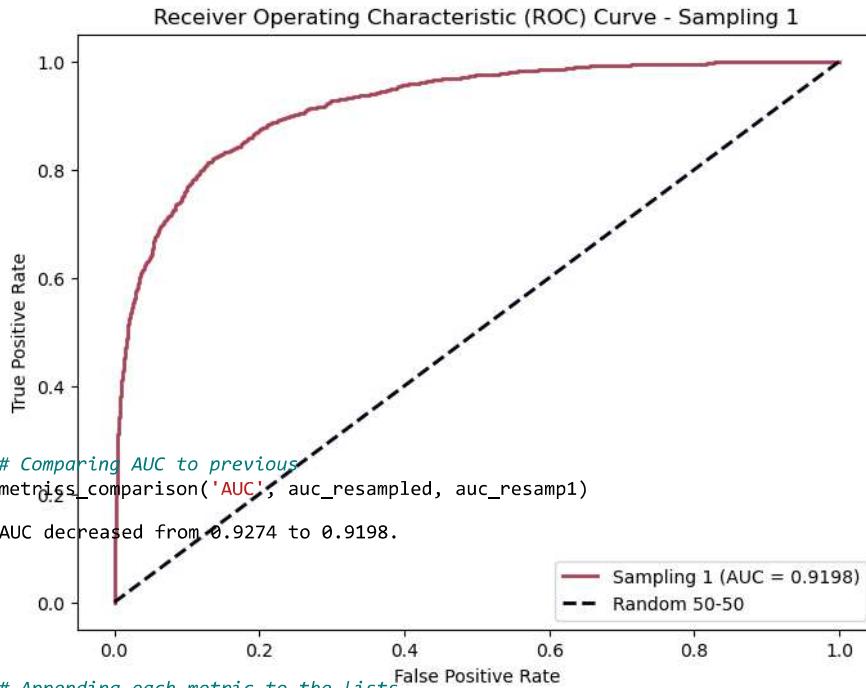


```
In [179]: # Calling confusion matrix for resampled model
rs1_cfn_matrix = confusion_matrix_display(rs1_pipeline, y_test, rs1_y_pred)
```



### 5) ROC Curve and AUC

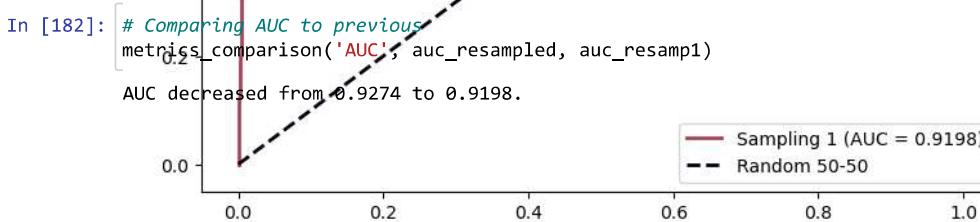
```
In [180]: # Plotting the ROC curve for the baseline model
resamp1_y_pred_proba = rs1_pipeline.predict_proba(X_test['original_review'])[:, 1]
auc_resamp1 = plot_roc_curve(y_test, resamp1_y_pred_proba, model_name=resampled1)
```



```
In [182]: # Comparing AUC to previous metrics
metrics_comparison('AUC', auc_resampled, auc_resamp1)
AUC decreased from 0.9274 to 0.9198.
```

```
In [183]: # Appending each metric to the lists
modelnames.append(resampled1_model_name)
accuracies.append(accuracy_rs1)
f1s.append(f1_rs1)
aucs.append(auc_rs1)
recalls.append(recall_rs1)
recall0_detractors.append(detectorive1power in distinguishing between
                           negative and positive classes.)
roc_aucs.append(auc_resamp1)
```

### 4th iteration: including stopwords



In [183]: # Appending each metric to the lists modelnames.append(resampled1\_model\_name)  
accuracys.append(accuracy\_rs1)  
f1s.append(f1\_rs1)  
aucs.append(auc\_resamp1)  
recalls.append(recall\_rs1)  
roc\_aucs.append(auc\_resamp1)

# Printing the interpretation of each metric to the user.

Recall\_0\_9274 indicates the model has good detection power in distinguishing between negative and positive reviews.

#### 4th iteration: including stopwords

I will now test fitting the vectorizer by removing the stopwords from reviews to see if this can help predictions be more accurate.

##### 1) Fitting and training train data

In [184]: # Defining the pipeline steps, including the stopwords\_list created  
tfidf\_vectorizer = TfidfVectorizer(stop\_words=stopwords\_list)  
naive\_bayes\_classifier = MultinomialNB()  
  
# Instantiating the pipeline with the undersampler and the new vectorizer  
pipeline\_nostop = imblearn.pipeline.Pipeline([  
 ('tfidf', tfidf\_vectorizer),  
 ('us', RandomUnderSampler(sampling\_strategy=1, random\_state=42)),  
 ('classifier', naive\_bayes\_classifier)  
])  
  
# Fit the pipeline on X\_train['original\_review'] and y\_train  
pipeline\_nostop.fit(X\_train['original\_review'], y\_train)  
  
# Generating predictions  
nostop\_y\_pred = pipeline\_nostop.predict(X\_test['original\_review'])

##### 2) Evaluation Metrics

In [185]: # Naming the model  
nostop\_model\_name = 'No Stopwords'  
  
# Calling the function and recording into the defined values  
accuracy\_nostop, f1\_nostop, recall\_nostop, cv\_nostop = evaluation\_metrics(  
 y\_test,  
 nostop\_y\_pred)  
print("Removing English stopwords improved all scores overall. I will now verify  
this by class.")  
X\_train['original\_review'].improved all scores overall. I will now verify how t  
true this is by class.  
  
Accuracy: 0.8319  
F1-Score: 0.8382  
Recall: 0.8319  
Mean Cross-Validated Accuracy: 0.8466

In [186]: # Comparing the classification metrics between baseline and Resampled model  
metrics\_comparison('AUC', auc\_resampled, auc\_nostop)  
metrics\_comparison('F1', f1\_rs1, f1\_nostop)  
metrics\_comparison('Recall', recall\_rs1, recall\_nostop)  
metrics\_comparison('Precision', precision\_rs1, precision\_nostop)  
metrics\_comparison('Cross-validated Accuracy', cv\_rs1, cv\_nostop)

Macro avg	0.78596	0.83238	0.80128	5123
Micro avg	0.8319	0.8319	0.8319	3773
F1	0.8319	0.8319	0.8319	1350
Precision	0.8319	0.8319	0.8319	3773
Recall	0.8319	0.8319	0.8319	1350
Support	5123	3773	3773	3773

Accuracy improved from 0.8319 to 0.8466.

```
In [187]: # Printing the overall summary
print('Removing English stopwords improved all scores overall. I will now verify
      pipeline_nostop,
      Removing English stopwords improved all scores overall. I will now verify how t
      rue this is by class.

Accuracy: 0.8319
F1-Score: 0.8382
Recall: Classification Report
Mean Cross-Validated Accuracy: 0.8466
```

```
In [188]: # Comparing the classification metrics between baseline and Resampled model
metrics_comparison('Precision', precision_rs1, precision_nostop)
metrics_comparison('Recall', recall_rs1, recall_nostop)
metrics_comparison('F1', f1_rs1, f1_nostop)
metrics_comparison('Cross-validated Accuracy', cv_rs1, cv_nostop)
```

Not Detractors improved from 0.8312 to 0.8319.	3773
F1 improved from 0.8312 to 0.8382.	1350
Recall improved from 0.8312 to 0.8319.	
Cross-Validated Accuracy improved from 0.8312 to 0.8466.	
macro avg	0.78596 0.83238 0.80128 5123
weighted avg	0.85554 0.83193 0.83820 5123

#### 4) Confusion Matrix

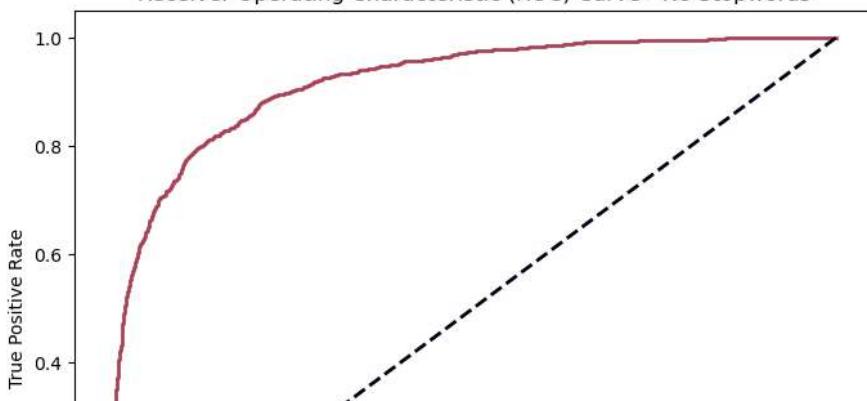
```
In [189]: # Calling the confusion matrix function
nostop_cfn_matrix = confusion_matrix_display(pipeline_nostop, y_test, nostop_y_pr
```

Model Performance: Confusion Matrix

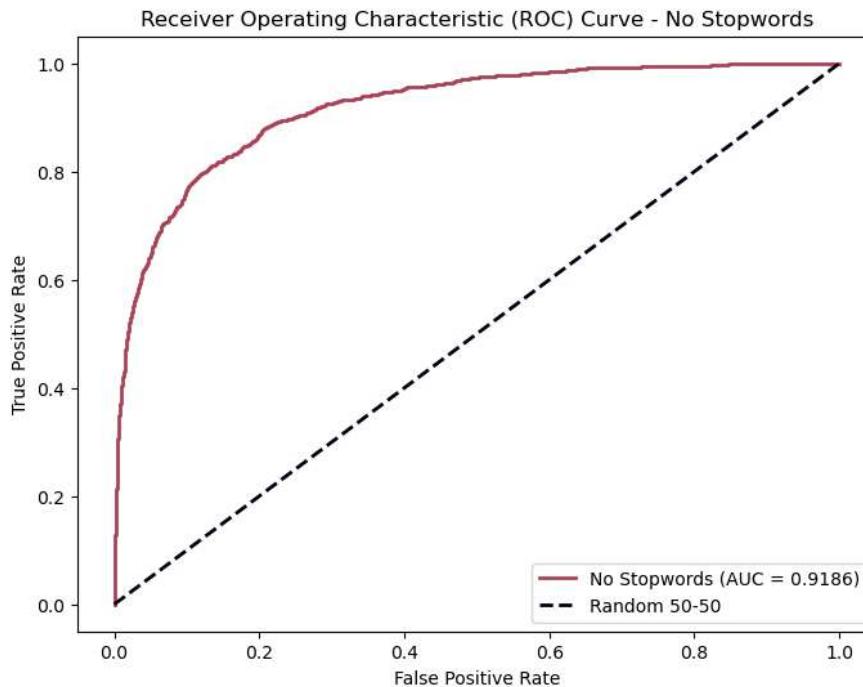


```
In [190]: # Plotting the ROC curve for the baseline model
nostop_y_pred_proba = pipeline_nostop.predict_proba(X_test['original_review'])[:, 1]
auc_nostop = plot_roc_curve(y_test, nostop_y_pred_proba, model_name=nostop_model)
Recall for detractors slightly decreased, but remains very high while all other metrics also
```

Receiver Operating Characteristic (ROC) Curve - No Stopwords



```
In [190]: # Plotting the ROC curve for the baseline model
nostop_y_pred_proba = pipeline_nostop.predict_proba(X_test['original_review'])[:, 1]
auc_nostop = plot_roc_curve(y_test, nostop_y_pred_proba, model_name=nostop_model)
Recall for detractors slightly decreased, but remains very high while all other metrics also
```



```
In [191]: # Printing the interpretation
auc_interpretation(auc_nostop)
```

AUC = 0.9186. The model has a great predictive power in distinguishing between the positive and negative classes.

```
In [192]: metrics_comparison('AUC', auc_resamp1, auc_nostop)
```

AUC decreased from 0.9198 to 0.9186.

```
In [193]: # Appending each metric to the lists
modelnames.append(nostop_model_name)
accuracies.append(accuracy_nostop)
f1s.append(f1_nostop)
recalls.append(recall_nostop)
recalls_detractors.append(recall_detractors_nostop)
cv_accs.append(cv_nostop)
roc_aucs.append(auc_nostop)
```

## 5th iteration: including only the industry specific stopwords

```
In [195]: # Defining the pipeline steps, including the stopwords list created
tfidf_vectorizer = TfidfVectorizer(stop_words=new_stopwords)
naive_bayes_classifier = MultinomialNB()
```

```
# This section adds the pipeline with the undersampling and the new vectorizer
pipeline_nonewstop = imblearn.pipeline.Pipeline([
    ('tfidf', tfidf_vectorizer),
    ('rus', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', naive_bayes_classifier)
])
```

```
In [194]: # Removing pipeline history of industry specific stopwords and y_train
pipeline_nonewstop.fit(X_train['original_review'], y_train)
```

```
Out[194]: ['hotel', 'room', 'night', 'day', 'stay', 'resort', 'place']
nonewstop_y_pred = pipeline_nonewstop.predict(X_test['original_review'])
```

### 2) Evaluation Metrics

```
In [195]: # For analysis purposes, I had removed a list of industry specific stopwords which came up often
# defining the pipeline steps, including the stopwords list created
tfidf_vectorizer = TfidfVectorizer(stop_words=new_stopwords)
naive_bayes_classifier = MultinomialNB()

# will test our model on by including only this list as stopwords
# including the original stopwords, and the new vectorizer
pipeline_nonewstop = imblearn.pipeline.Pipeline([
    ('tfidf', tfidf_vectorizer),
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', naive_bayes_classifier)
])

In [194]: # Reminding the list of industry specific stopwords and y_train
# Pipeline with newstop.fit(X_train['original_review'], y_train)

Out[194]: ['hotel', 'room', 'night', 'day', 'stay', 'resort', 'place']
nonewstop_y_pred = pipeline_nonewstop.predict(X_test['original_review'])
```

## 2) Evaluation Metrics

```
In [196]: # Naming the model
nonewstop_model_name = 'Industry Stopwords'

# Calling the function and recording into the defined values
accuracy_nonewstop, f1_nonewstop, recall_nonewstop, cv_nonewstop = evaluation_me
    y_test,
    nonewstop_y_pred,
    pipeline_nonewstop,
    X_train['original_review'],
    y_train)

Accuracy: 0.8331
F1-Score: 0.8396
Recall: 0.8331
Mean Cross-Validated Accuracy: 0.8460
```

```
In [197]: # Comparing the evaluation metrics between baseline and Resampled model
metrics_comparison('Accuracy', accuracy_nostop, accuracy_nonewstop)
metrics_comparison('F1', f1_nostop, f1_nonewstop)
metrics_comparison('Recall', recall_nostop, recall_nonewstop)
metrics_comparison('Cross-validated Accuracy', cv_nostop, cv_nonewstop)
```

Accuracy improved from 0.8319 to 0.8331.  
 F1 improved from 0.8382 to 0.8396.  
 Recall improved from 0.8319 to 0.8331.  
 Cross-validated Accuracy decreased from 0.8466 to 0.8460.

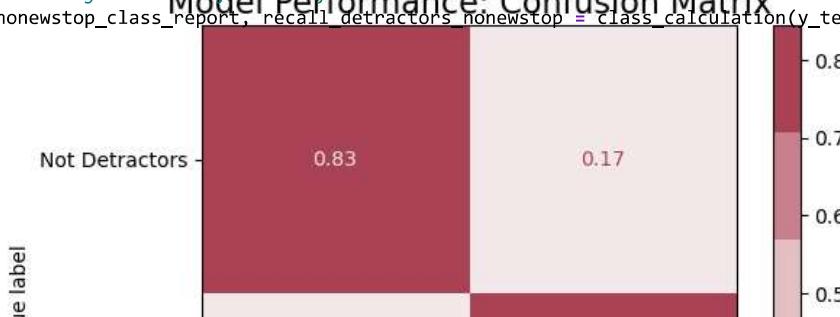
All scores continued increasing, except from cross-validated accuracy. Cross-validated accuracy is an estimate of the model's performance on test data using cross-validation. The decrease is very small, but it indicates the model would perform slightly lower on unseen data.

```
In [198]: # Printing the overall summary
print('Including Industry stopwords improved scores on training data compared to')

Including Industry stopwords improved scores on training data compared to incl
        4) Confusion Matrix
```

```
In [200]: # 3) Classification Report
# Calling the confusion matrix function
nonewstop_cfn_matrix = confusion_matrix_display(pipeline_nonewstop, y_test, none
```

```
In [199]: # Calling the classification function
nonewstop_class_report, recall_detractors_nonewstop = class_calculation(y_test, r
```

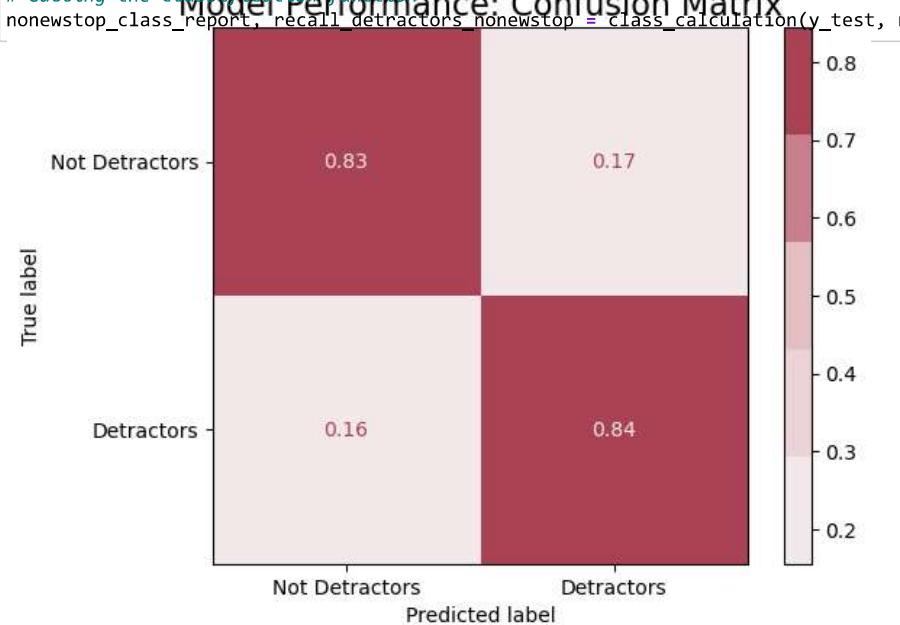


Including Industry stopwords improved scores on training data compared to including English stopwords.

In [200]: # Calling the confusion matrix function

```
nonewstop_cfn_matrix = confusion_matrix_display(pipeline_nonewstop, y_test, nonewstop)
```

In [199]: # Calling the classification function



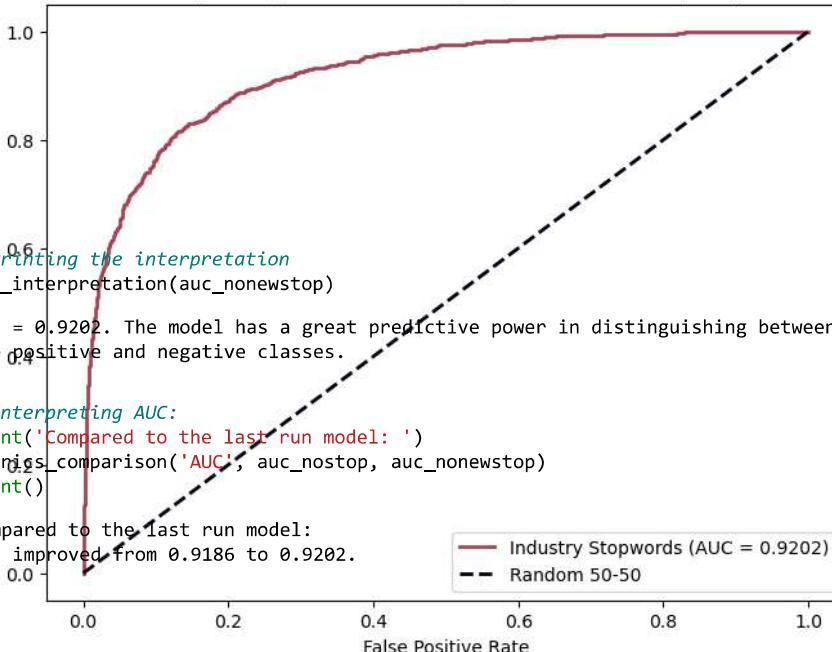
Removing stopwords increased the recall for Detractors back to higher results, above 84%.

## 5) ROC Curve and AUC

In [201]: # Plotting the ROC curve for the baseline model

```
nonewstop_y_pred_proba = pipeline_nonewstop.predict_proba(X_test['original_review'])
auc_nonewstop = plot_roc_curve(y_test, nonewstop_y_pred_proba, model_name=nonewstop)
```

Receiver Operating Characteristic (ROC) Curve - Industry Stopwords



In [202]: # Printing the interpretation

```
auc_interpretation(auc_nonewstop)
```

AUC = 0.9202. The model has a great predictive power in distinguishing between the positive and negative classes.

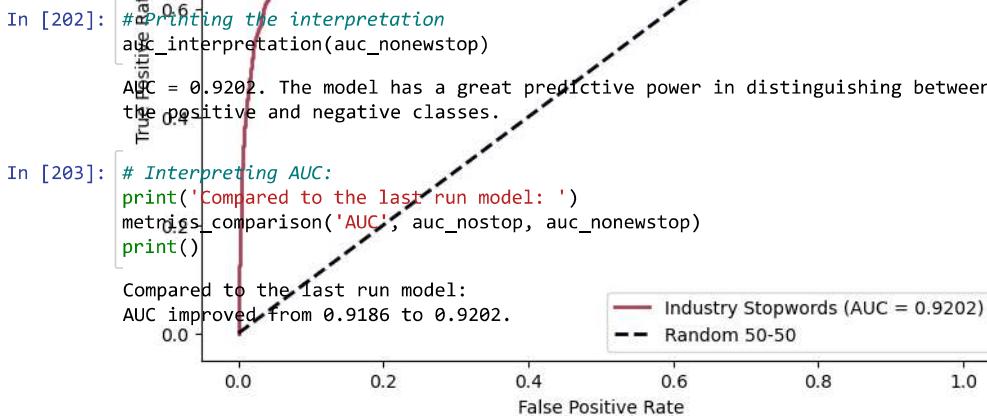
In [203]: # Interpreting AUC:

```
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_nostop, auc_nonewstop)
print()
```

Compared to the last run model:  
AUC improved from 0.9186 to 0.9202.

In [204]: # Appending each metric to the lists

```
modelnames.append(nonewstop_model_name)
accuracies.append(accuracy_nonewstop)
```



In [204]: `# Appending each metric to the lists  
modelnames.append(nonewstop_model_name)  
accuracies.append(accuracy_nonewstop)  
f1s.append(f1_nonewstop)  
recalls.append(recall_nonewstop)  
recalls_detractors.append(recall_detractors_nonewstop)  
cv_accs.append(cv_nonewstop)  
roc_aucs.append(auc_nonewstop)`

## 6th iteration: adding both English and industry stopwords

I will now test the model performance after excluding the combined list of stopwords: English, and industry-related.

### 1) Fitting and training train data

I will create a new list of stopwords as a copy of the existing one, so the industry stopwords can be added to it.

In [205]: `# Reminding the List of English stopwords  
stopwords_list[:3] # is the one including all English stopwords.`

Out[205]: `['i', 'me', 'my']`

In [206]: `# Reminding the List of industry specific stopwords  
new_stopwords # includes industry specific stopwords`

Out[206]: `['hotel', 'room', 'night', 'day', 'stay', 'resort', 'place']`

In [207]: `# Creating a new list to combine both stopwords lists  
full_stopwords = stopwords_list.copy()`

In [209]: `# Adding the new_stopwords to the existing list  
full_stopwords.extend(new_stopwords) # using the new_stopwords_list created`

In [208]: `tfidf_vectorizer = TfidfVectorizer(stop_words=full_stopwords)  
naive_bayes_classifier = MultinomialNB()`

`# Verifying the list of new_stopwords were added to the stopwords list  
assert full_stopwords[-len(new_stopwords)] == new_stopwords`

`# Instantiating the pipeline with the undersampler and the new vectorizer  
pipeline_fullstop = imblearn.pipeline.Pipeline([  
 ('tfidf', tfidf_vectorizer),  
 ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),  
 ('classifier', naive_bayes_classifier)  
)`

`# Fit the pipeline on X_train['original_review'] and y_train  
pipeline_fullstop.fit(X_train['original_review'], y_train)`

`# Generating predictions  
fullstop_y_pred = pipeline_fullstop.predict(X_test['original_review'])`

### 2) Evaluation Metrics

```
# Adding the new_stopwords to the existing list
In [209]: full_stopwords.append(new_stopwords) # Using the new_stopwords_list created
tfidf_vectorizer = TfidfVectorizer(stop_words=full_stopwords)

# Verifying the list of new_stopwords were added to the stopwords list
In [208]: assert full_stopwords[-len(new_stopwords):] == new_stopwords
# Instantiating the pipeline with the undersampler and the new vectorizer
pipeline_fullstop = imblearn.pipeline.Pipeline([
    ('tfidf', tfidf_vectorizer),
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', naive_bayes_classifier)
])

# Fit the pipeline on X_train['original_review'] and y_train
pipeline_fullstop.fit(X_train['original_review'], y_train)

# Generating predictions
fullstop_y_pred = pipeline_fullstop.predict(X_test['original_review'])
```

## 2) Evaluation Metrics

```
In [210]: # Naming the model
fullstop_model_name = 'Full Stopwords'

# Calling the function and recording into the defined values
accuracy_fullstop, f1_fullstop, recall_fullstop, cv_fullstop = evaluation_metrics(
    y_test,
    fullstop_y_pred,
    pipeline_fullstop,
    X_train['original_review'],
    y_train)
```

Accuracy: 0.8345  
F1-Score: 0.8405  
Recall: 0.8345  
Mean Cross-Validated Accuracy: 0.8481

The accuracy score now increased slightly, however it remains below just guessing the majority class.

```
In [211]: # Comparing the evaluation metrics between baseline and Resampled model
metrics_comparison('Accuracy', accuracy_nonewstop, accuracy_fullstop)
metrics_comparison('F1', f1_nonewstop, f1_fullstop)
metrics_comparison('Recall', recall_nonewstop, recall_fullstop)
metrics_comparison('Cross-validated Accuracy', cv_nonewstop, cv_fullstop)
```

Accuracy improved from 0.8331 to 0.8345.  
F1 improved from 0.8396 to 0.8405.  
Recall improved from 0.8331 to 0.8345.  
Cross-validated Accuracy improved from 0.8460 to 0.8481.

```
In [212]: # Printing the overall summary
print('Including full list of stopwords improved scores.')
```

Including full list of stopwords improved scores.

- 4) Confusion Matrix
- 3) Classification Report

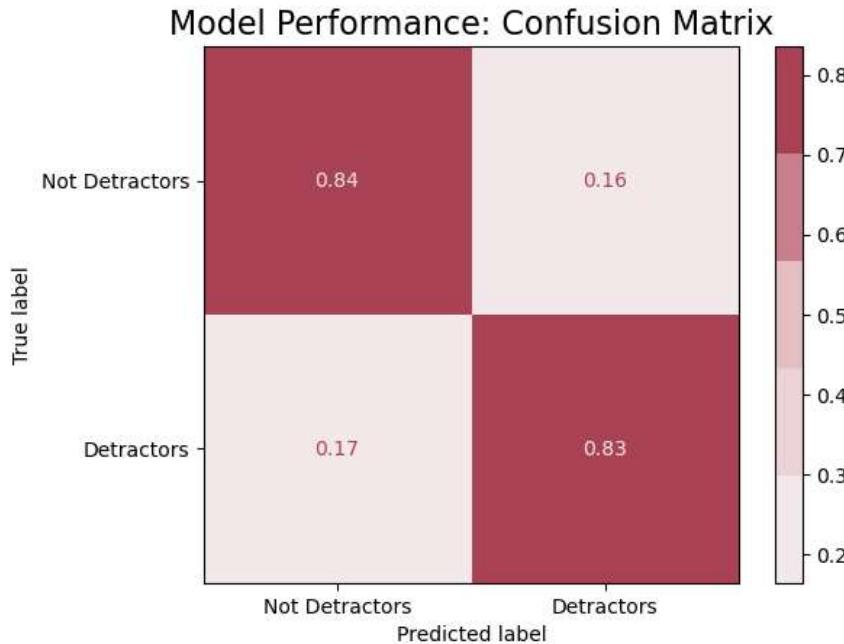
```
In [214]: # Calling the confusion_matrix function
In [213]: fullstop_cfmatrix.confusion_matrix_display(pipeline_fullstop, y_test, fullstop_y_pred, fullstop_class_report, recall_detractors_fullstop = class_calculation(y_test, fullstop_y_pred))
```

Model Performance: Confusion Matrix



- 4) Confusion Matrix  
3) Classification Report

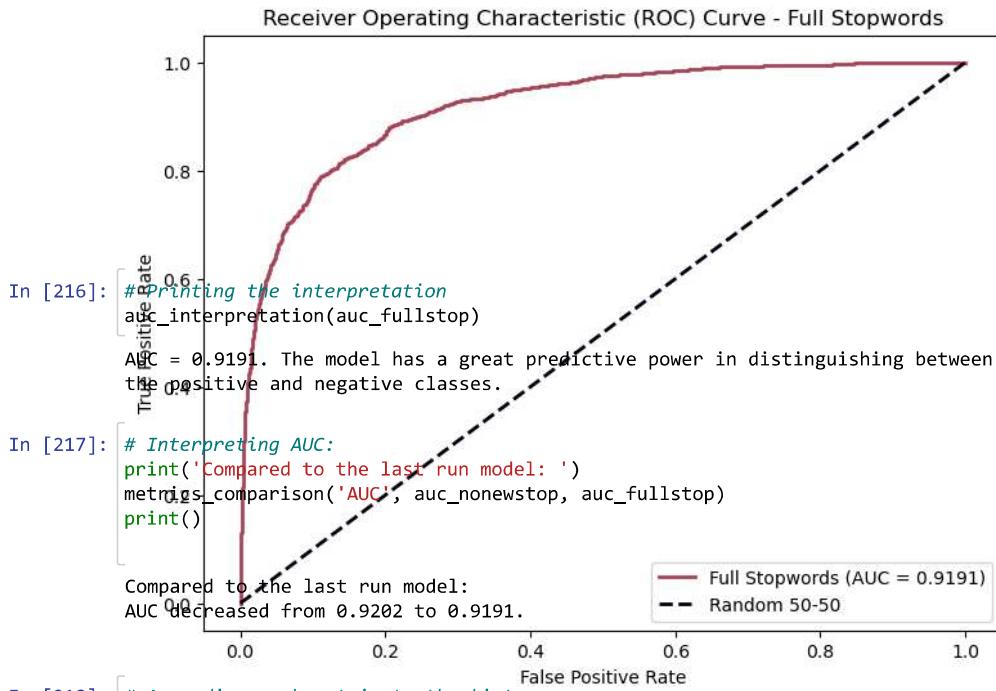
```
In [214]: # Calling the confusion_matrix function
In [213]: fullstop_confusion_matrix_display(pipeline_fullstop, y_test, fullstop_class_report, recall_detractors_fullstop = class_calculation(y_test, fu
```



Including the full list improved the model but not recall for detractors.

- 5) ROC Curve and AUC

```
In [215]: # Plotting the ROC curve for the baseline model
fullstop_y_pred_proba = pipeline_fullstop.predict_proba(X_test['original_review'])
auc_fullstop = plot_roc_curve(y_test, fullstop_y_pred_proba, model_name=fullstop
```



```
In [216]: # Printing the interpretation
auc_interpretation(auc_fullstop)
AUC = 0.9191. The model has a great predictive power in distinguishing between the positive and negative classes.

In [217]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_nonewstop, auc_fullstop)
print()

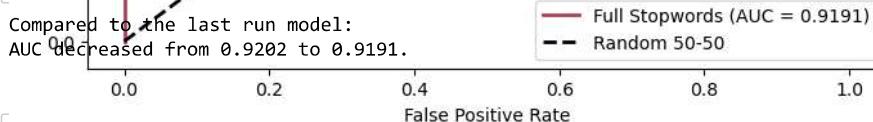
Compared to the last run model:
AUC decreased from 0.9202 to 0.9191.

In [218]: # Appending each metric to the lists
modelnames.append(fullstop_model_name)
accuracies.append(accuracy_fullstop)
f1s.append(f1_fullstop)
recalls.append(recall_fullstop)
```

```
In [216]: # Printing the interpretation
auc_interpretation(auc_fullstop)

AUC = 0.9191. The model has a great predictive power in distinguishing between the positive and negative classes.
```

```
In [217]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_nonewstop, auc_fullstop)
print()
```



```
In [218]: # Appending each metric to the lists
modelnames.append(fullstop_model_name)
accuracies.append(accuracy_fullstop)
f1s.append(f1_fullstop)
recalls.append(recall_fullstop)
recalls_detractors.append(recall_detractors_fullstop)
cv_accs.append(cv_fullstop)
roc_aucs.append(auc_fullstop)
```

## 7th iteration: Lemmatized Reviews

Models perform better when industry stopwords are excluded from the reviews. Let's now review how it compares when it runs of lemmatized words.

1) Fitting and training train data

For the exploratory analysis, some adjectives were included in this list, however they would be important to help predicting the sentiment of the review. They will be dropped from this list so they are still considered by the model.

```
In [219]: # Only lemmatizing X_train['lemm_review']
X_train['lemmatized_review'] = X_train['original_review'].apply(lambda x: tokenized(x))
X_train['lemmatized_review'] = X_train['lemmatized_review'].apply(lambda tokens:
```

```
In [220]: # Storing outside a list
X_train['lemmatized_review'] = X_train['lemmatized_review'].apply(lambda x: ' '.join(x))
```

```
In [221]: # Viewing the newly created column
X_train[:3]
```

```
Out[221]:
Review Rating label original_review Review_prep_nolist lemmatized_review
-----
```

```
In [222]: # Applying the same to X_test
X_test['lemmatized_review'] = X_test['original_review'].apply(lambda x: tokenized(x))
X_test['lemmatized_review'] = X_test['lemmatized_review'].apply(lambda tokens: ' '.join(tokens))
X_test['breakfast_good_review'] = X_test['breakfast_good_review'].apply(lambda x: ' '.join(x))
```

```
In [223]: # Defining the pipeline steps, including stop words and no stop words
tfidf_vectorizer = TfidfVectorizer(stop_words='english')
naive_bayes_classifier = MultinomialNB()
# Instantiating the pipeline with the undersampler and the new vectorizer
```

```
pipeline_lemm = imblearn.pipeline.Pipeline([
    ('tfidf', tfidf_vectorizer),
    ('undersampler', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', naive_bayes_classifier)
])
```

```
# Fit the pipeline on X_train['original_review'] and y_train
pipeline_lemm.fit(X_train['lemmatized_review'], y_train)
```

```
In [222]: # Applying the same to X_test
X_test['original_review'] = X_test['original_review'].apply(lambda x: tokenizer(x))
X_test['lemmatized_review'] = X_test['lemmatized_review'].apply(lambda x: lemmatizer(x))
X_test['breakfast_good_review'] = X_test['breakfast_good_review'].apply(lambda x: lemmatizer(x))

In [223]: # Defining the pipeline steps, including the words that are lemmatized
tfidf_vectorizer = TfidfVectorizer(stop_words='english')
naive_bayes_classifier = MultinomialNB()
# Instantiating the pipeline with the undersampler and the new vectorizer
pipeline_lemm = imblearn.pipeline.Pipeline([
    ('tfidf', tfidf_vectorizer),
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', naive_bayes_classifier)
])

# Fit the pipeline on X_train['original_review'] and y_train
pipeline_lemm.fit(X_train['lemmatized_review'], y_train)

# Generating predictions
lemm_y_pred = pipeline_lemm.predict(X_test['lemmatized_review'])
```

## 2) Evaluation Metrics

```
In [224]: # Naming the model
lemm_model_name = 'Lemmatized'

# Calling the function and recording into the defined values
accuracy_lemm, f1_lemm, recall_lemm, cv_lemm = evaluation_metrics(
    y_test,
    lemm_y_pred,
    pipeline_lemm,
    X_train['lemmatized_review'],
    y_train)

Accuracy: 0.8271
F1-Score: 0.8344
Recall: 0.8271
Mean Cross-Validated Accuracy: 0.8382
```

```
In [225]: # Comparing the evaluation metrics between undersampled model (best model) and no
metrics_comparison('Accuracy', accuracy_fullstop, accuracy_lemm)
metrics_comparison('F1', f1_fullstop, f1_lemm)
metrics_comparison('Recall', recall_fullstop, recall_lemm)
metrics_comparison('Cross-validated Accuracy', cv_fullstop, cv_lemm)

Accuracy decreased from 0.8345 to 0.8271.
F1 decreased from 0.8405 to 0.8344.
Recall decreased from 0.8345 to 0.8271.
Cross-validated Accuracy decreased from 0.8481 to 0.8382.
```

```
In [226]: # Printing the overall summary
In [227]: print(classification_report(y_test, lemm_y_pred))
lemm_class_report, recall_detractors_lemm = class_calculation(y_test, lemm_y_pred)
Most scores decreased slightly.
Classification Report:
precision    recall    f1-score   support
Not Detractors  0.94076  0.81659  0.87429      3773
          Detractors  0.62554  0.85630  0.72295      1350

           accuracy          0.82705      5123
          macro avg   0.78315   0.83644   0.79862      5123
     weighted avg   0.85770   0.82705   0.83441      5123
```

At first glance, it would seem like the model performance decreased. It actually got the highest results ever recorded for the Detractors class, which is what we are aiming for.

## 4) Confusion Matrix

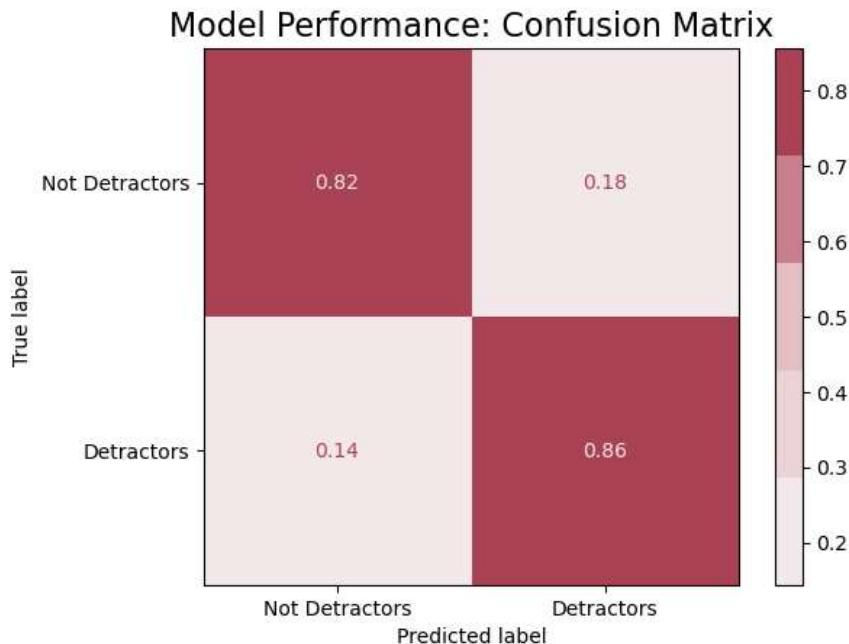
```
In [226]: # Printing the overall summary
In [227]: # Printing most the classification function.
lemm_class_report, recall_detractors_lemm = class_calculation(y_test, lemm_y_pred)
Most scores decreased slightly.
Classification Report:
precision    recall    f1-score   support
3) Classification Report
Not Detractors      0.94076   0.81659   0.87429   3773
Detractors          0.62554   0.85630   0.72295   1350

accuracy                   0.82705   5123
macro avg                 0.78315   0.83644   0.79862   5123
weighted avg              0.85770   0.82705   0.83441   5123
```

At first glance, it would seem like the model performance decreased. It actually got the highest results ever recorded for the Detractors class, which is what we are aiming for.

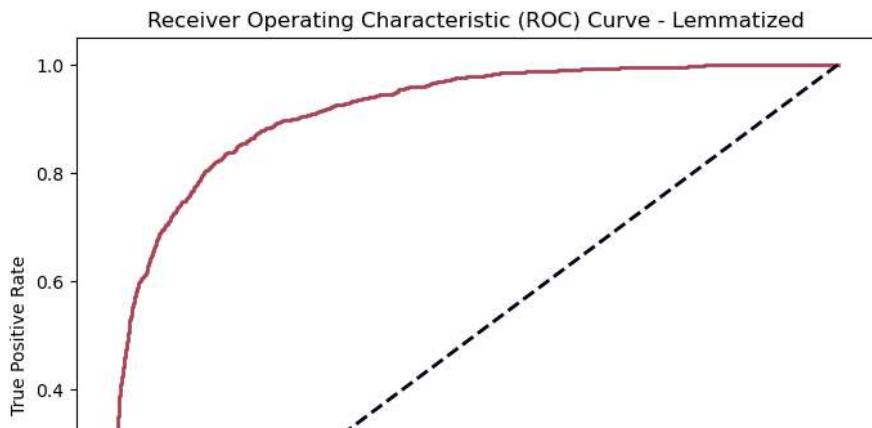
#### 4) Confusion Matrix

```
In [228]: # Calling the confusion matrix function
lemm_cfn_matrix = confusion_matrix_display(pipeline_lemm, y_test, lemm_y_pred)
```

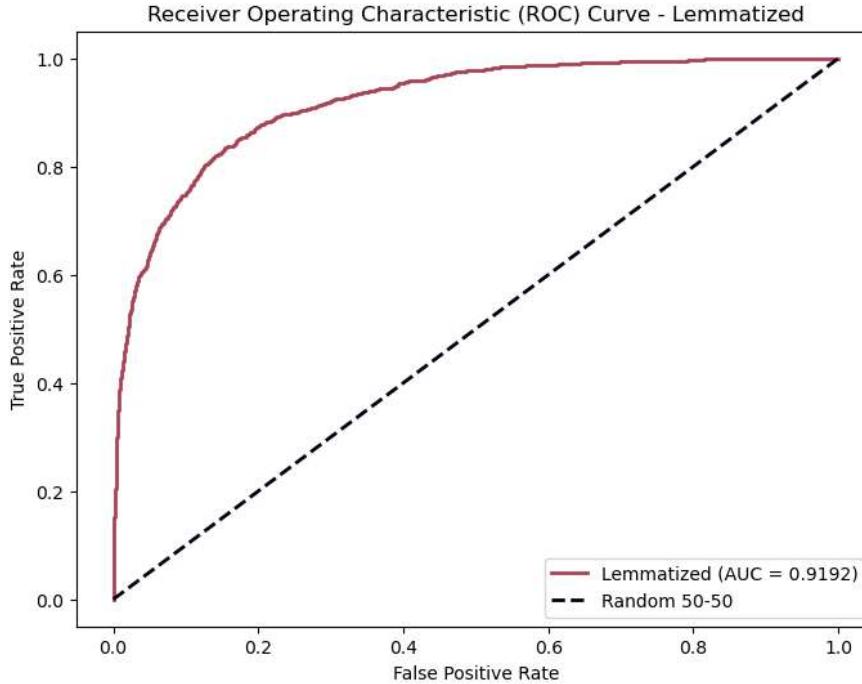


#### 5) ROC Curve and AUC

```
In [229]: # Plotting the ROC curve for the baseline model
lemm_y_pred_proba = pipeline_lemm.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_lemm = plot_roc_curve(y_test, lemm_y_pred_proba, model_name=lemm_model_name)
```



```
In [229]: # Plotting the ROC curve for the baseline model
lemm_y_pred_proba = pipeline_lemm.predict_proba(X_test['lemmatized_review'])[:, :]
auc_lemm = plot_roc_curve(y_test, lemm_y_pred_proba, model_name=lemm_model_name)
```



```
In [230]: # Printing the interpretation
auc_interpretation(auc_lemm)
```

AUC = 0.9192. The model has a great predictive power in distinguishing between the positive and negative classes.

```
In [231]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_fullstop, auc_lemm)
print()
```

Compared to the last run model:  
AUC improved from 0.9191 to 0.9192.

```
In [232]: # Appending each metric to the lists
modelnames.append(lemm_model_name)
accuracies.append(accuracy_lemm)
f1s.append(f1_lemm)
recalls.append(recall_lemm)
recalls_detractors.append(recall_detractors_lemm)
cv_accs.append(cv_lemm)
roc_aucs.append(auc_lemm)
```

```
In [234]: # Defining the pipeline steps
Seeing the total results for the Detractors class, it makes sense that the AUC improved compared to the previous classifier = MultinomialNB()
```

```
# Instantiating the pipeline with the undersampler and the new vectorizer
pipeline_prep = imblearn.pipeline.Pipeline([
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', naive_bayes_classifier)
])
```

1) Fitting and training on train data

```
# Fitting the pipeline on X_train and X_test['Review_prep_noList'] which contains
pipeline_prep.fit(X_train['Review_prep_noList'], y_train)
```

```
# Not applying stopwords to the vectorizer, as they were applied to the preprocessed reviews
```

```
# Generating predictions
prep_y_pred = pipeline_prep.predict(X_test['Review_prep_noList'])
```

## 2) Evaluation Metrics

```
In [234]: # Defining the pipeline steps
Since the total results for detractors class, it makes sense that the AUC improved compared
to the previous classifier = MultinomialNB()

# Instantiating the pipeline with the undersampler and the new vectorizer
pipeline_prep = imblearn.pipeline.Pipeline([
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', naive_bayes_classifier)
])
1) Fitting and training on train data

# Fitting the pipeline on X_train and X_test['Review_prep_nolist'] which contains
# Not applying stopwords to the vectorizer, as they were applied to the preprocessed
In [233]: # Generating predictions
prep_y_pred = pipeline_prep.predict(X_test['Review_prep_nolist'])
```

## 2) Evaluation Metrics

```
In [235]: pp_model_name = 'Preprocessed'

# Calling the function and recording into the defined values
accuracy_pp, f1_pp, recall_pp, cv_pp = evaluation_metrics(
    y_test,
    prep_y_pred,
    pipeline_prep,
    X_train['Review_prep_nolist'],
    y_train)

Accuracy: 0.8271
F1-Score: 0.8341
Recall: 0.8271
Mean Cross-Validated Accuracy: 0.8404
```

```
In [236]: # Comparing the evaluation metrics between undersampled model (best model) and no
metrics_comparison('Accuracy', accuracy_rs1, accuracy_pp)
metrics_comparison('F1', f1_rs1, f1_pp)
metrics_comparison('Recall', recall_rs1, recall_pp)
metrics_comparison('Cross-validated Accuracy', cv_rs1, cv_pp)

Accuracy decreased from 0.8312 to 0.8271.
F1 decreased from 0.8378 to 0.8341.
Recall decreased from 0.8312 to 0.8271.
Cross-validated Accuracy decreased from 0.8442 to 0.8404.
```

```
In [237]: # Printing the overall summary
print('All scores decreased on the fully preprocessed reviews.')

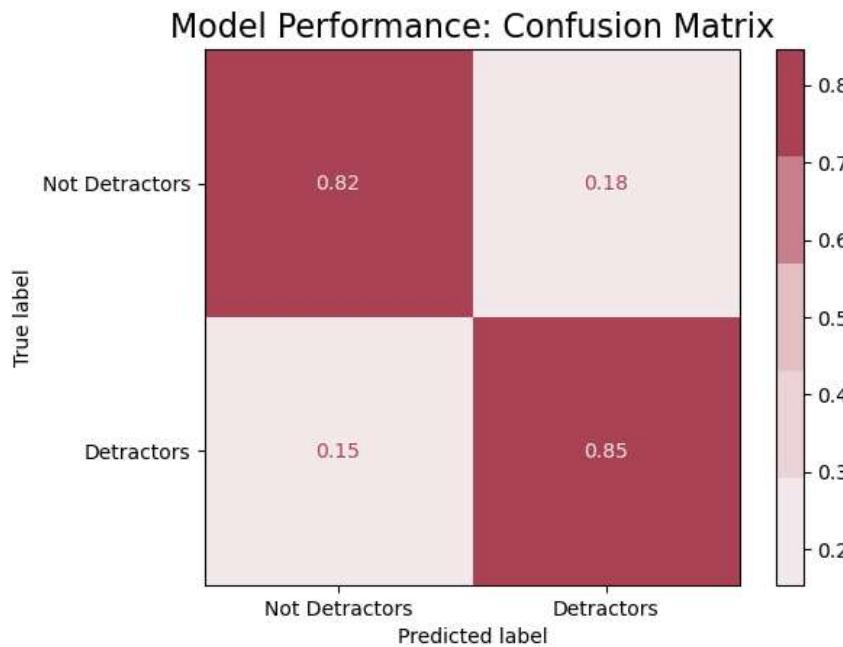
All scores decreased on the fully preprocessed reviews.
```

## 3) Classification Report

```
In [238]: # Blotting the confusion matrix report
prep_cmate_report, confusion_matrix_display(pipeline_prep, X_test, prep_y_pred)
```



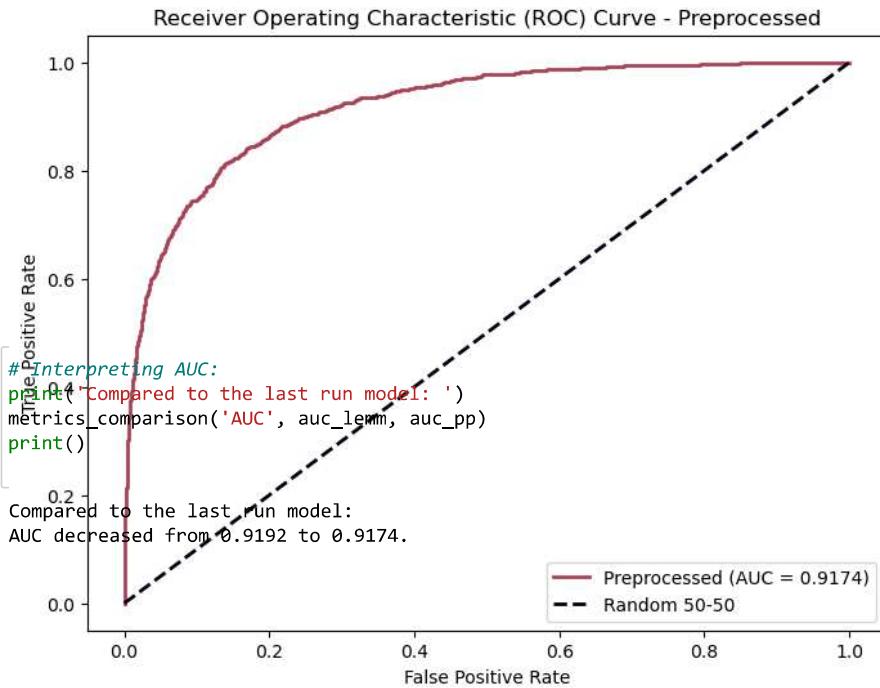
```
In [238]: # Drawing the confusion matrix report
prep_cv_report, confusion_report_display(pipeline_prep_cv, test_prep, prep_cv)
```



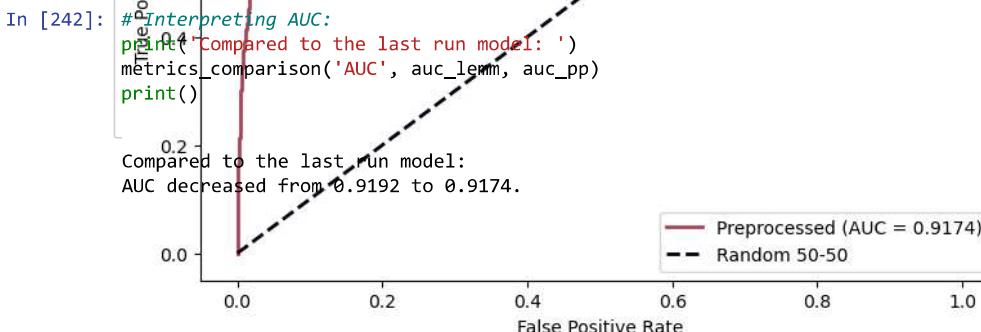
This is also reflected on recall for detractors which decreased compared to the one for lemmatized reviews.

### 5) ROC Curve and AUC

```
In [240]: # Plotting the ROC curve for the baseline model
prep_y_pred_proba = pipeline_prep.predict_proba(X_test['Review_prep_nolist'])[:, 1]
auc_pp = plot_roc_curve(y_test, prep_y_pred_proba, model_name=pp_model_name)
```



```
In [243]: # Appending each metric to the lists
modelnames.append(pp_model_name)
accs.append(cv_pp)
f1s.append(f1_pp)
AUCs.append(AUC_pp)
fpr.append(fpr_pp)
tprs.append(tprs_pp)
cv_acccs.append(cv_pp)
```



In [243]: # Appending each metric to the lists

In [241]:

```
modelnames.append(pp_modelname)
ACC.append(accuracy_pp)
f1s.append(f1_pp)
AUC.append(auc_pp) # AUC of 0.9174 means the model has a great predictive power in distinguishing between the positive and negative reviews
cv_accs.append(cv_pp)
roc_aucs.append(auc_pp)
```

## 9th iteration: Tuning Tfifd Vectorizer - Hyperparameter tuning

The model performed better when stopwords were removed but worse when applied on the full tokenized reviews. Let's try to use combinatoric grid searching to find the best parameters for the vectorizer.

### 1) Fitting and training on train data

In [244]:

```
from sklearn.model_selection import GridSearchCV
```

```
# Defining the pipeline steps, excluding any manual input of features, with the values defined earlier
# The pipeline still includes the undersampler to ensure class imbalance is addressed
gs_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', MultinomialNB())
])

# Define the parameter grid to search over
parameters = {
    'tfidf_max_features': [100, None],
    'tfidf_max_df': [0.7, 0.9],
    'tfidf_min_df': [1, 3],
    'tfidf_sublinear_tf': [True, False],
    # 'tfidf_ngram_range': [(1, 1), (2, 2)],
    # including the list of stopwords that was defined earlier
    'tfidf_stop_words': [stopwords_list, new_stopwords, full_stopwords],
}
```

In [246]: # Recording the best parameters and printing them

In [245]:

```
best_tfidf_params = grid_search.best_params_
grid_search = GridSearchCV(gs_pipeline, parameters, cv=5, scoring='recall', error_score=0)
# Recording the best estimator as the best_pipeline
best_pipeline = grid_search.best_estimator_
# fitting the pipeline on the best training data
grid_search.fit(X_train['lemmatized_review'], y_train)
print("Best Parameters: ", best_tfidf_params)
```

Out[245]:

```
Best Parameters: {'tfidf_max_df': 0.7, 'tfidf_max_features': None, 'tfidf_min_df': 1, 'tfidf_stop_words': ['hotel', 'room', 'night', 'day', 'stay', 'resource'], 'estimator': Pipeline(sublinear_tf': True}
    ↗ TfidfVectorizer
```

In [247]: # Fitting the best pipeline on training data

```
best_pipeline.fit(X_train['lemmatized_review'], y_train)
```

```
# Fitting the best model on the full training data
```

```
best_gs_y_pred = best_pipeline.predict(X_test['lemmatized_review'])
```

### 2) Evaluation Metrics

```
In [246]: # Recording the best parameters and printing them
In [245]: best_fitting_params = grid_search.best_params_
grid_search = GridSearchCV(gs_pipeline, parameters, cv=5, scoring='recall', error_score=0)
# Recording the best estimator as the best_pipeline
best_pipeline = grid_search.best_estimator_
grid_search.fit(X_train['lemmatized_review'], y_train)
print("Best Parameters: ", best_fitting_params)
```

```
Out[245]: Best Parameters: {'tfidf__max_df': 0.7, 'tfidf__max_features': None, 'tfidf__min_df': 1, 'tfidf__stop_words': ['hotel', 'room', 'night', 'day', 'stay', 'resource'], 'estimator': Pipeline(sublinear_tf': True)}
```

```
In [247]: # Fitting the best pipeline on training data
best_pipeline.fit(X_train['lemmatized_review'], y_train)

# Fitting the best model on the full training data
best_gs_y_pred = best_pipeline.predict(X_test['lemmatized_review'])
```

## 2) Evaluation Metrics

```
In [248]: # Naming the model
gs_model_name = 'NaiveBayes GS'

# Calling the function and recording into the defined values
accuracy_gs, f1_gs, recall_gs, cv_gs = evaluation_metrics(
    y_test,
    best_gs_y_pred,
    best_pipeline,
    X_train['lemmatized_review'],
    y_train)
```

Accuracy: 0.8274  
F1-Score: 0.8349  
Recall: 0.8274  
Mean Cross-Validated Accuracy: 0.8370

```
In [249]: # Comparing the evaluation metrics between undersampled model (best model) and no
metrics_comparison('Accuracy', accuracy_fullstop, accuracy_gs)
metrics_comparison('F1', f1_fullstop, f1_gs)
metrics_comparison('Recall', recall_fullstop, recall_gs)
metrics_comparison('Cross-validated Accuracy', cv_fullstop, cv_gs)
```

Accuracy decreased from 0.8345 to 0.8274.  
F1 decreased from 0.8405 to 0.8349.  
Recall decreased from 0.8345 to 0.8274.  
Cross-validated Accuracy decreased from 0.8481 to 0.8370.

```
In [250]: # Printing the overall summary
print('Tuning hyperparameters for recall did not increase overall results.')
```

Tuning hyperparameters for recall did not increase overall results.

## 3) Classification Report 4) Confusion Matrix

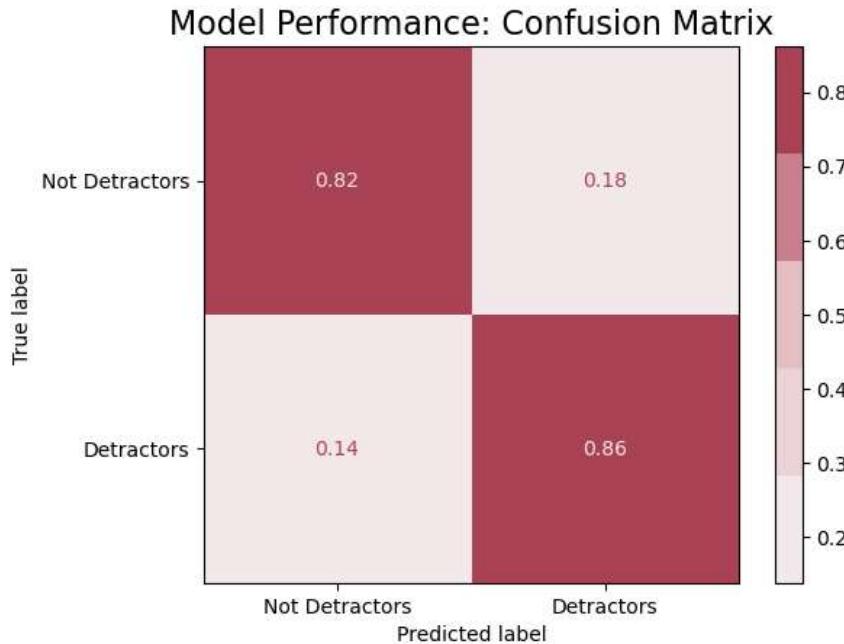
```
In [251]: # Calling the classification report function
In [252]: # Fitting the report function
gs_cfn_matrix = confusion_matrix_display(best_pipeline, y_test, best_gs_y_pred)
Classification Report:
```

Model Performance: Confusion Matrix



- 3) Classification Report  
4) Confusion Matrix

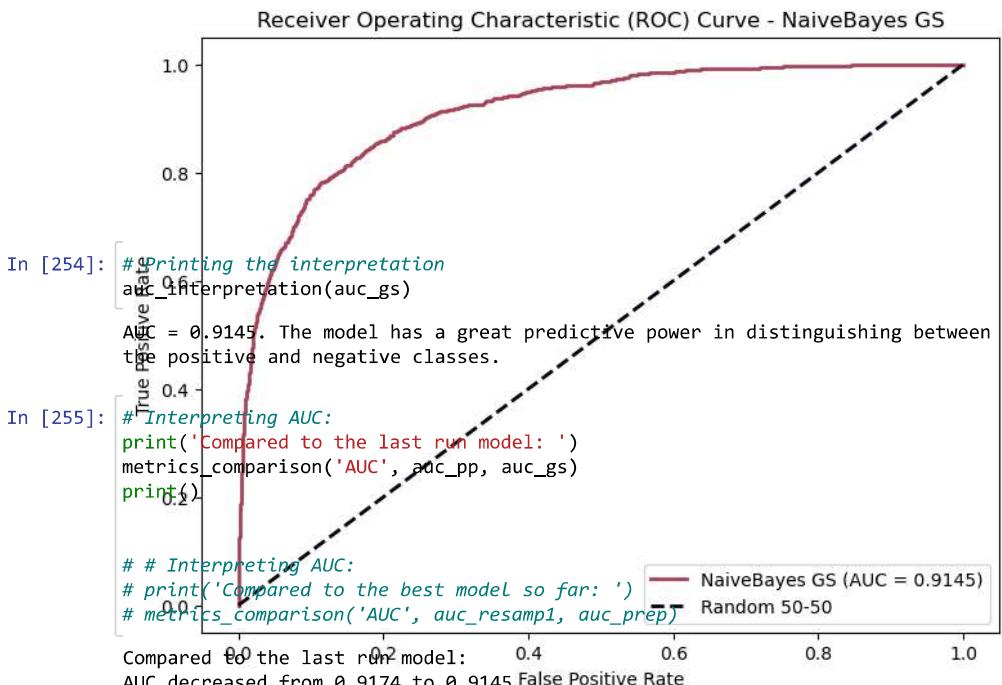
```
In [251]: # Calling the classification report function
In [252]: # Plotting the confusion matrix
gs_cfn_matrix = confusion_matrix_display(best_pipeline, y_test, best_gs_y_pred)
Classification Report:
```



Recall for detractors was the highest ever reached despite lower overall results. This is the best model so far.

- 5) ROC Curve and AUC

```
In [253]: # Plotting the ROC curve for the baseline model
gs_y_pred_proba = best_pipeline.predict_proba(X_test['original_review'])[:, 1]
auc_gs = plot_roc_curve(y_test, gs_y_pred_proba, model_name=gs_model_name)
```



```
In [254]: # Printing the interpretation
auc_gs_interpretation(auc_gs)

AUC = 0.9145. The model has a great predictive power in distinguishing between the positive and negative classes.
```

```
In [255]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_pp, auc_gs)
print()

# # Interpreting AUC:
# print('Compared to the best model so far: ')
# metrics_comparison('AUC', auc_resamp1, auc_prep)
```

Every time the true positive rate decreases for not detractors AUC decreases slightly as well

```
In [254]: # Printing the interpretation
auc_gs_interpretation(auc_gs)

AUC = 0.9145. The model has a great predictive power in distinguishing between
the positive and negative classes.
```

```
In [255]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_pp, auc_gs)
print()

# # Interpreting AUC:
# print('Compared to the best model so far: ')
# metrics_comparison('AUC', auc_resamp1, auc_prep)
```

Compared to the last run model: AUC decreased from 0.9174 to 0.9145

Every time the true positive rate decreases for not\_detractors, AUC decreases slightly as well.

```
In [256]: # Appending each metric to the lists
modelnames.append(gs_model_name)
accuracies.append(accuracy_gs)
f1s.append(f1_gs)
recalls.append(recall_gs)
recalls_detractors.append(recall_detractors_gs)
cv_acccs.append(cv_gs)
roc_aucs.append(auc_gs)
```

## 5. d) TfifdVectorizer and K-Nearest Neighbor

### 10th iteration: K-Nearest Neighbor Tfifd Vectorizer

#### 1) Fitting and training on train data

The previous model was a bit computationally expensive. Let's see if the simpler K-Nearest Neighbor classifier would improve on that end. Nevertheless, kNN makes predictions based on what similar cases around it suggest so there is a risk it captures more noise created by the imbalanced dataset, despite the undersampled not\_detractors reviews.

Since the grid search model with Naive Bayes recorded the best results, the iterations will be made on this base.

For higher computing performance, the best parameters recorded on the vectorizer with Multinomial Naive Bayes will be kept. Only the classifier will be modified. Let's see if, by using the best TFIDF parameters with another classifier, I can improve further these predictions.

```
In [258]: # Creating a function that modifies the parameters, to allow them to be used for
# Verifying the best vectorizer parameters
def transform_params(best_params):
    print(best_params)
    new_best_tfidf_params = {}
    for key, value in best_params.items():
        if 'tfidf_max_df' in key:
            # Removing 'tfidf' from the key
            new_key = key.replace('tfidf__', '_')
            new_best_params[new_key] = value
    return new_best_params
```

The parameters were transformed to be used in the vectorizer but they are not needed at this stage.

```
In [259]: # Getting the newly defined parameters
new_best_tfidf_params = transform_params(best_tfidf_params)

# # Inspecting them
# print(new_best_tfidf_params)
```

```
In [260]: from sklearn.neighbors import KNeighborsClassifier

# Defining the pipeline with the fixed tfidf parameters and RandomForestClassifier
knn_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', TfidfVectorizer(stop_words=new_stopwords)),
    ('knn', KNeighborsClassifier(n_neighbors=5))])
```

```
In [258]: # Creating a function that modifies the parameters, to allow them to be used for
In [257]: def transform_params(best_params):
    print(best_params)
    new_best_params = {}
    for key, value in best_params.items():
        if key == 'tfidf_max_df': None, 'tfidf_min_df': 1, 'tfidf_
        stop_words': ['hotel', 'room', 'night', 'day', 'stay', 'resort', 'place'], 'tf
        idf_sublinear_tf': True}
        new_best_params[new_key] = value
    return new_best_params
```

The parameters were transformed to be used in the vectorizer but they are not needed at this

```
In [259]: # Setting the newly defined parameters
# new_best_tfidf_params = transform_params(best_tfidf_params)

# # Inspecting them
# print(new_best_tfidf_params)
```

```
In [260]: from sklearn.neighbors import KNeighborsClassifier

# Defining the pipeline with the fixed tfidf parameters and RandomForestClassifier
knn_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', TfidfVectorizer(stop_words=new_stopwords)),
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', KNeighborsClassifier())
])

# Fitting the pipeline on training data
knn_pipeline.fit(X_train['lemmatized_review'], y_train)

# Making predictions on test data
knn_y_pred = knn_pipeline.predict(X_test['lemmatized_review'])
```

## 2) Evaluation Metrics

```
In [261]: # Naming the model and calling the function to evaluate it
knn_model_name = 'kNN'

# Calling the function and recording into the defined values
accuracy_knn, f1_knn, recall_knn, cv_knn = evaluation_metrics(
    y_test,
    knn_y_pred,
    knn_pipeline,
    X_train['lemmatized_review'],
    y_train
)
```

Accuracy: 0.7451  
F1-Score: 0.7569  
Recall: 0.7451  
Mean Cross-Validated Accuracy: 0.7564

```
In [262]: # Comparing the evaluation metrics between undersampled model (best model) and no
metrics_comparison('Accuracy', accuracy_gs, accuracy_knn)
```

```
In [264]: # Calling the function
metrics_comparison('F1', f1_gs, f1_knn)
metrics_comparison('recall', recall_gs, recall_knn)
knn_classification_report(y_test, knn_y_pred)
metrics_comparison('Cross-Validated Accuracy', cv_gs, cv_knn)
```

Classification Report:  
Accuracy decreased from 0.8274 to 0.7451.  
F1 decreased from 0.8329 to 0.7569.  
recall decreased from 0.8274 to 0.7451.  
Not Detractors 0.88535 0.75113 from 0.81273 to 0.3773.  
Cross-Validated Accuracy decreased from 0.8370 to 0.7564.  
Detractors 0.51145 0.72815 0.60086 1350

```
In [263]: # Printing the overall summary
print("The metrics recorded from kNN are so far behind what was recorded, no iteration
      weighted avg 0.78682 0.74507 0.75690 5123
The metrics recorded from kNN are so far behind what was recorded, no iteration
      s will be built from this model.
```

## 4) Confusion Matrix 3) Classification Report

```
In [265]: # Recording and displaying the confusion matrix
knn_confusion_matrix = confusion_matrix_display(knn_pipeline, y_test, knn_y_pred)
```

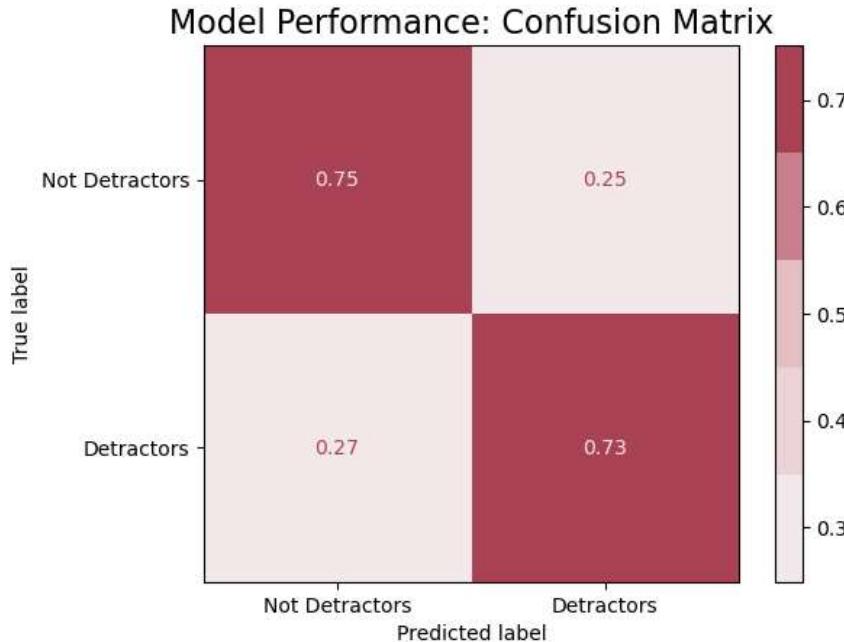
```
In [264]: # Calling the function
metrics_comparison('F1', f1_gs, f1_knn)
knn_class_report(recall_gs, recall_knn)
metrics_comparison('Cross-Validated Accuracy', ev_gs, cv_knn_y_pred)

Classification Report:
Accuracy decreased from 0.8274 to 0.7451.
F1 decreased from 0.8349 to 0.7569.
Precision 0.8349  f1-score   support
recall decreased from 0.8274 to 0.7451.
Not Detractors 0.88535 0.75113 0.81273
Cross-Validated Accuracy decreased from 0.8370 to 0.7564.
Detractors 0.51145 0.72815 0.60086 1350
```

```
In [263]: # Printing the overall summary
print(metrics_overview(gs, knn))
The metrics recorded from kNN are so far behind what was recorded, no iterations will be built from this model.
```

- 4) Confusion Matrix
- 3) Classification Report

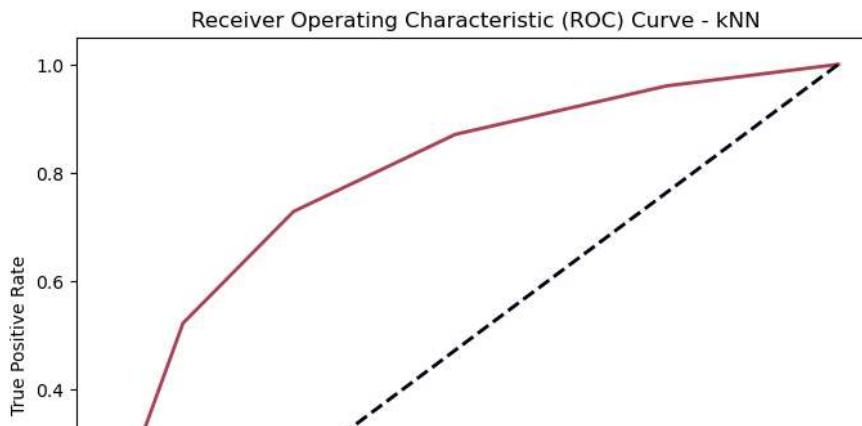
```
In [265]: # Recording and displaying the confusion matrix
knn_confusion_matrix = confusion_matrix_display(knn_pipeline, y_test, knn_y_pred)
```



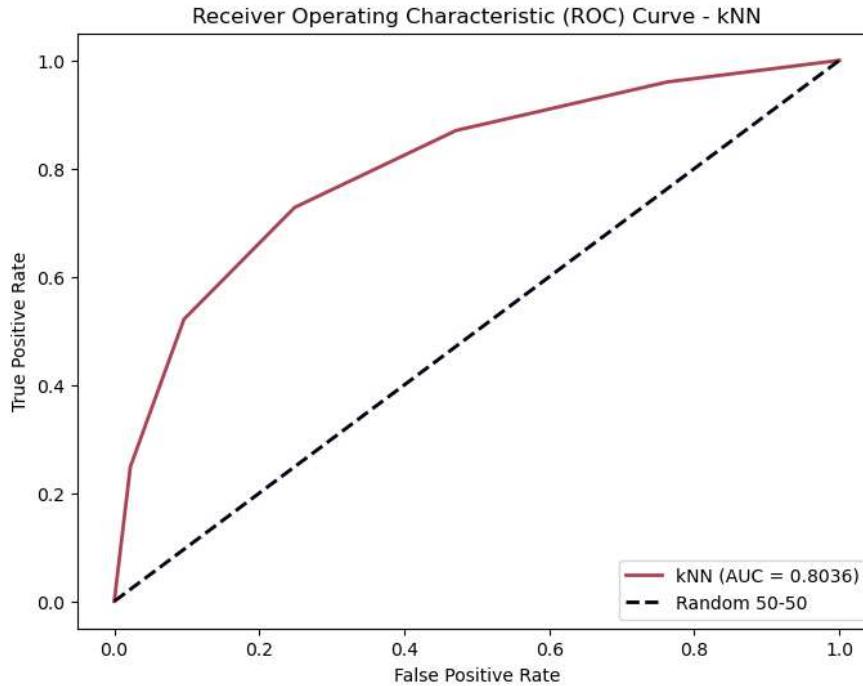
K-Nearest Neighbor classifier had a negative impact on the model. There is no point in trying to tune it.

- 5) ROC Curve and AUC

```
In [266]: # Plotting the ROC curve for the baseline model
knn_y_pred_proba = knn_pipeline.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_knn = plot_roc_curve(y_test, knn_y_pred_proba, model_name=knn_model_name)
```



```
In [266]: # Plotting the ROC curve for the baseline model
knn_y_pred_proba = knn_pipeline.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_knn = plot_roc_curve(y_test, knn_y_pred_proba, model_name=knn_model_name)
```



```
In [267]: # Printing the interpretation
auc_interpretation(auc_knn)
```

AUC = 0.8036. The model has a good predictive power.

```
In [268]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_gs, auc_knn)
print()
```

Compared to the last run model:  
AUC decreased from 0.9145 to 0.8036.

```
In [269]: # Appending each metric to the lists
modelnames.append(knn_model_name)
accuracies.append(accuracy_knn)
f1s.append(f1_knn)
recalls.append(recall_knn)
recalls_detractors.append(recall_detractors_knn)
cv_accs.append(cv_knn)
roc_aucs.append(auc_knn)
```

```
In [270]: # Importing the relevant packages
from sklearn.tree import DecisionTreeClassifier
5. c) TfIdfVectorizer and Decision Trees
# Defining the pipeline with new classifier, but the same best parameters
dt_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', TfIdfVectorizer(stop_words='english')),
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', DecisionTreeClassifier())
])
Decision trees work well for understanding language because they are easy to interpret and handle the nuances in how words relate. They are good at understanding what words matter most and can deal with different types of word data without much difficulty.
# Fitting the pipeline on training data
dt_pipeline.fit(X_train['lemmatized_review'], y_train)

# Making predictions on test data
dt_y_pred = dt_pipeline.predict(X_test['lemmatized_review'])
```

## 2) Evaluation Metrics

```
In [271]: # Naming the model
```

```
In [270]: # Importing the relevant packages
from sklearn.tree import DecisionTreeClassifier
5. c) TfIdfVectorizer and Decision Trees
# Defining the pipeline with new classifier, but the same best parameters
dt_pipeline = imblearn.pipeline.Pipeline([
    11th iteration: Decision Trees and TfIdfVectorizer,
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', DecisionTreeClassifier())
])
Decision trees work well for understanding language because they are easy to interpret and handle the nuances in how words relate. They are good at understanding what words matter most and can deal with different types of word data without much difficulty.
dt_pipeline.fit(X_train['lemmatized_review'], y_train)

# Making predictions on test data
dt_y_pred = dt_pipeline.predict(X_test['lemmatized_review'])
```

## 2) Evaluation Metrics

```
In [271]: # Naming the model
dt_model_name = 'DecisionTree'

# Calling the function and recording into the defined values
accuracy_dt, f1_dt, recall_dt, cv_dt = evaluation_metrics(
    y_test,
    dt_y_pred,
    dt_pipeline,
    X_train['lemmatized_review'],
    y_train)
```

Accuracy: 0.7293  
F1-Score: 0.7421  
Recall: 0.7293  
Mean Cross-Validated Accuracy: 0.7244

```
In [272]: # Comparing the evaluation metrics between undersampled model (best model) and no
metrics_comparison('Accuracy', accuracy_knn, accuracy_dt)
metrics_comparison('F1', f1_knn, f1_dt)
metrics_comparison('Recall', recall_knn, recall_dt)
metrics_comparison('Cross-validated Accuracy', cv_knn, cv_dt)
```

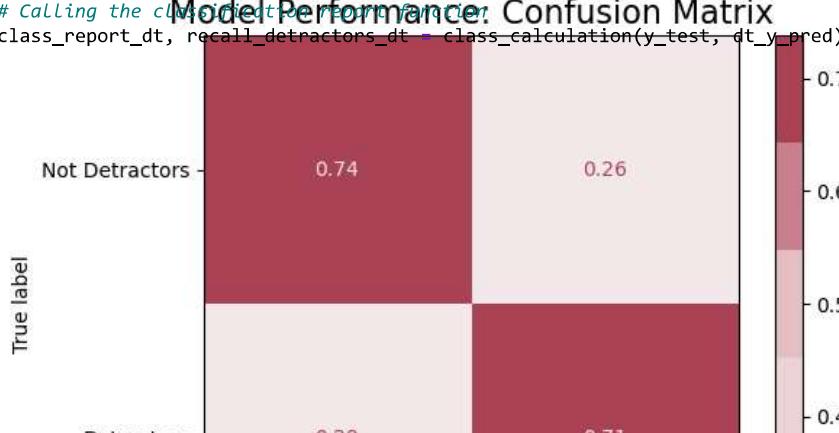
Accuracy decreased from 0.7451 to 0.7293.  
F1 decreased from 0.7569 to 0.7421.  
Recall decreased from 0.7451 to 0.7293.  
Cross-validated Accuracy decreased from 0.7564 to 0.7244.

```
In [273]: # Printing the overall summary
print('Decision tree produced results that were even lower than kNN. No iterations will be built from this.')
```

Decision tree produced results that were even lower than kNN. No iterations will be built from this.

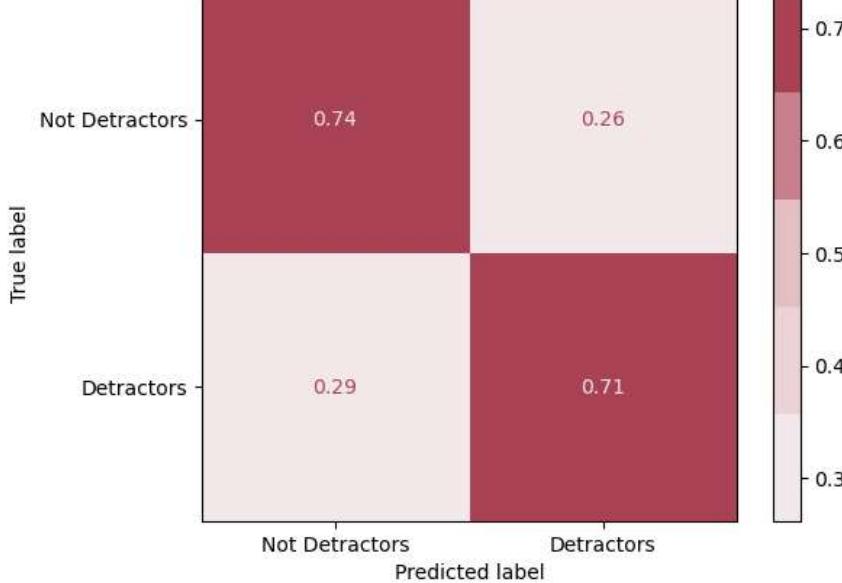
```
In [275]: # Recreating and displaying the confusion matrix
dt_cfn_matrix = confusion_matrix_display(dt_pipeline, y_test, dt_y_pred)
```

```
In [274]: # Calling the classification report function
class_report_dt, recall_detractors_dt = class_calculation(y_test, dt_y_pred)
```



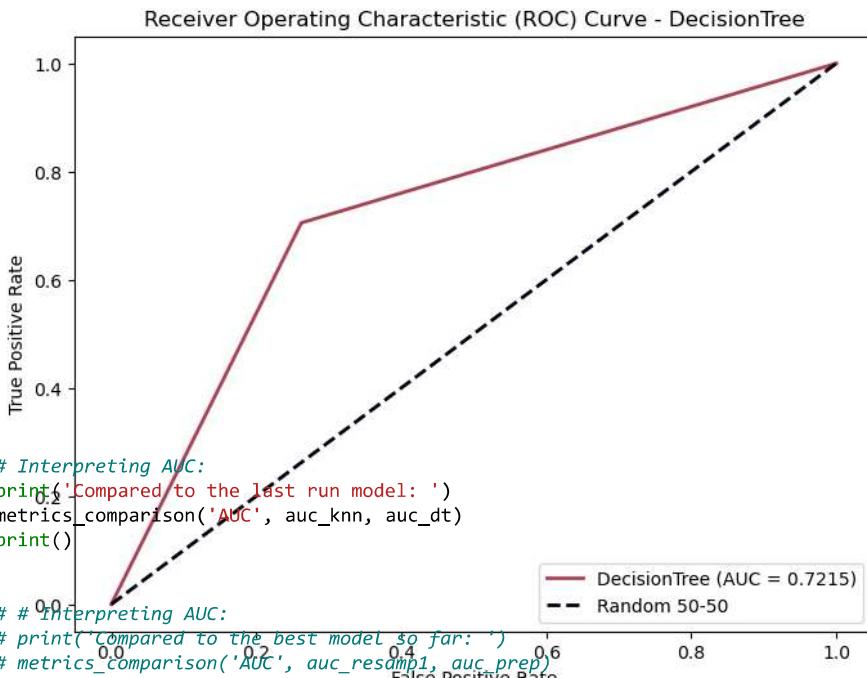
```
In [275]: # 3) Classification Report: the confusion matrix
dt_cfn_matrix = confusion_matrix_display(dt_pipeline, y_test, dt_y_pred)
```

```
In [274]: # Calling the classification performance function
class_report_dt, recall_detractors_dt = class_calculation(y_test, dt_y_pred)
```



## 5) ROC Curve and AUC

```
In [276]: # Plotting the ROC curve for the baseline model
dt_y_pred_proba = dt_pipeline.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_dt = plot_roc_curve(y_test, dt_y_pred_proba, model_name=dt_model_name)
```



```
In [278]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_knn, auc_dt)
print()
```

```
# # Interpreting AUC:
# print('Compared to the best model so far: ')
# metrics_comparison('AUC', auc_resamp1, auc_prep)
```

Compared to the last run model:

AUC decreased from 0.8036 to 0.7215.

```
In [277]: # Printing the interpretation
auc_interpretation(auc_dt)
```

```
In [279]: # AUC is 0.7215, so the model has some predictive power.
modelnames.append(dt_model_name)
accuracies.append(accuracy_dt)
f1s.append(f1_dt)
recalls.append(recall_dt)
recalls_detractors.append(recall_detractors_dt)
cv_accs.append(cv_dt)
```

```
In [278]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_knn, auc_dt)
print()

# # Interpreting AUC:
# print('Compared to the best model so far: ')
# metrics_comparison('AUC', auc_resamp1, auc_prep)
```

Compared to the last run model:  
AUC decreased from 0.8036 to 0.7215.

```
In [277]: # Printing the interpretation
auc_interpretation(auc_dt)
```

```
In [279]: # Appending each model has some predictive power.
modelnames.append(dt_model_name)
accuracies.append(accuracy_dt)
f1s.append(f1_dt)
recalls.append(recall_dt)
recalls_detractors.append(recall_detractors_dt)
cv_accs.append(cv_dt)
roc_aucs.append(auc_dt)
```

## 5. d) TfidfVectorizer and Random Forest

### 12th iteration: RandomForestClassifier Tuning Tfidf Vectorizer

#### 1) Fitting and training on train data

Random Forest classifiers can be thought of as an extension of multiple decision trees working together together to understand language text. Let's see if, by using the industry stopwords and lemmatized reviews with another classifier, I can improve further these predictions.

```
In [280]: # Importing the relevant package
from sklearn.ensemble import RandomForestClassifier

# Defining the pipeline with the fixed tfidf parameters and RandomForestClassifier
rf_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', TfidfVectorizer(stop_words=new_stopwords)),
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', RandomForestClassifier())
])

# Fitting the pipeline on training data
rf_pipeline.fit(X_train['lemmatized_review'], y_train)

# Making predictions on test data
rf_y_pred = rf_pipeline.predict(X_test['lemmatized_review'])
```

```
In [282]: # Comparing the evaluation metrics between the last model and the new one
metrics_comparison('Accuracy', accuracy_dt, accuracy_rf)
metrics_comparison('F1', f1_dt, f1_rf)
```

```
In [281]: metrics_comparison('Recall', recall_dt, recall_rf)
metrics_comparison('Cross-Validated Accuracy', cv_dt, cv_rf)
```

Accuracy improved from 0.7293 to 0.8366.

F1 improved from 0.7421 to 0.8421  
Recall improved from 0.7293 to 0.8366 ± evaluation\_metrics(  
Cross-Validated Accuracy improved from 0.7244 to 0.8507.

```
In [283]: # Printing pipeline overall summary
print(X_train['lemmatized_review'])
y_train
Random Forest finally started to improve results at similar levels than NaiveBayes.
yes.
Accuracy: 0.8366
F1-Score: 0.8421
Recall: 0.8366
3) Classification Report
Mean Cross-Validated Accuracy: 0.8507
```

```
In [282]: # Comparing the evaluation metrics between the Last model and the new one
metrics_comparison('Accuracy', accuracy_dt, accuracy_rf)
```

```
In [281]: metrics_comparison('Recall', recall_dt, recall_rf)
metrics_comparison('F1', f1_dt, f1_rf)
```

Accuracy improved from 0.7293 to 0.8366.

F1 improved from 0.7421 to 0.8421

Recall improved from 0.7293 to 0.8366 evaluation metrics

Cross-Validated Accuracy improved from 0.7244 to 0.8507.

```
In [283]: # Printing pipeline overall summary
```

```
print(rf_pipeline.named_steps['rf'].get_params())
y_train
```

Random Forest finally started to improve results at similar levels than NaiveBayes yes.

Accuracy: 0.8366

F1-Score: 0.8421

Recall: 0.8366

Mean Cross-Validated Accuracy: 0.8507

```
In [284]: # Calling the classification report function
```

```
rf_class_report, recall_detractors_rf = class_calculation(y_test, rf_y_pred)
```

Classification Report:

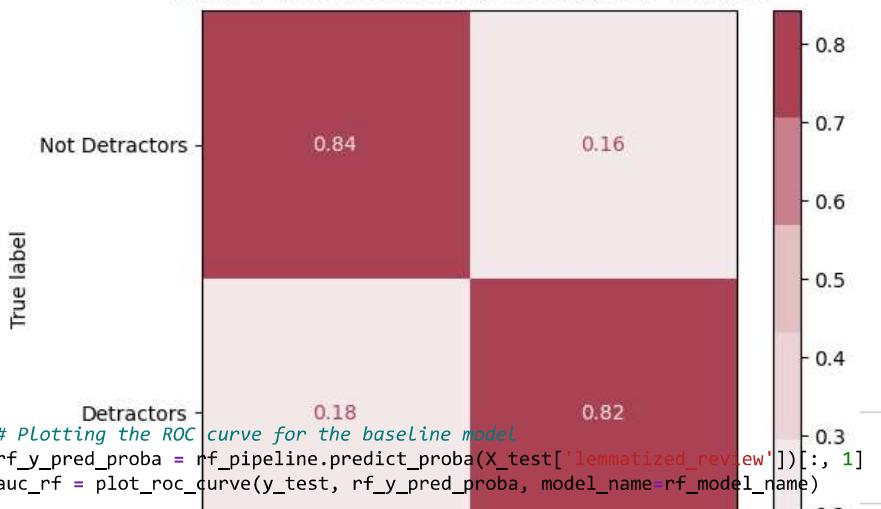
	precision	recall	f1-score	support
Not Detractors	0.92949	0.84204	0.88360	3773
Detractors	0.65044	0.82148	0.72602	1350
accuracy			0.83662	5123
macro avg	0.78997	0.83176	0.80481	5123
weighted avg	0.85596	0.83662	0.84208	5123

#### 4) Confusion Matrix

```
In [285]: # Recording and displaying the confusion matrix
```

```
rf_confusion_matrix = confusion_matrix_display(rf_pipeline, y_test, rf_y_pred)
```

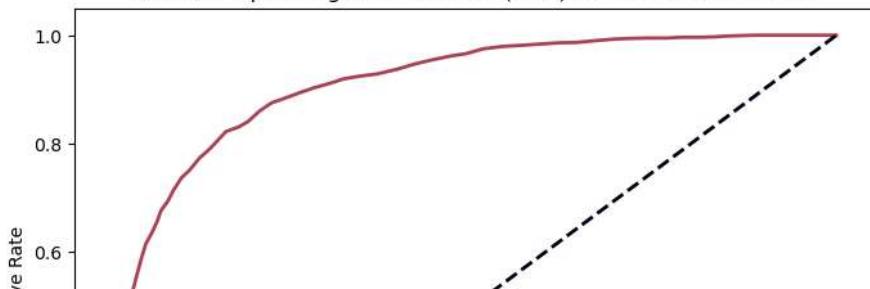
Model Performance: Confusion Matrix

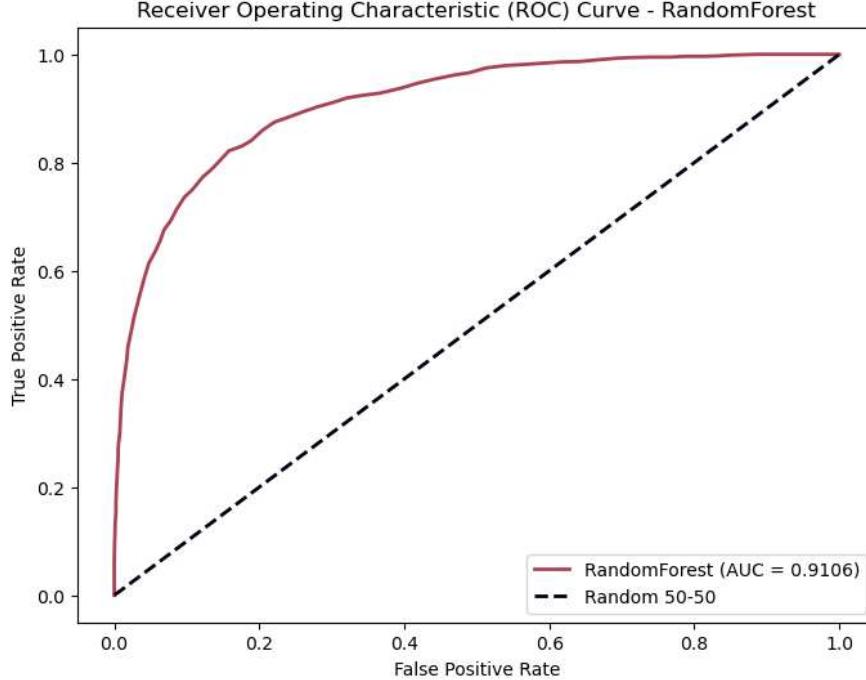
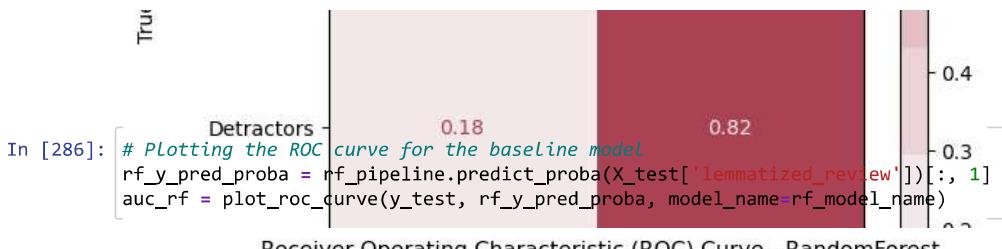


```
In [286]: # Plotting the ROC curve for the baseline model
```

```
rf_y_pred_proba = rf_pipeline.predict_proba(X_test[['lemmatized_review']])[:, 1]
auc_rf = plot_roc_curve(y_test, rf_y_pred_proba, model_name=rf_model_name)
```

Receiver Operating Characteristic (ROC) Curve - RandomForest





In [287]: # Printing the interpretation  
auc\_interpretation(auc\_rf)

AUC = 0.9106. The model has a great predictive power in distinguishing between the positive and negative classes.

In [288]: # Interpreting AUC:  
print('Compared to the last run model: ')  
metrics\_comparison('AUC', auc\_dt, auc\_rf)  
print()  
  
# # Interpreting AUC:  
# print('Compared to the best model so far: ')  
# metrics\_comparison('AUC', auc\_resamp1, auc\_prep)

Compared to the last run model:  
AUC improved from 0.7215 to 0.9106.

Gradient boosting is a class of ensemble algorithms constructed from decision tree models. Trees

In [289]: # Append current model to the lists  
model\_names.append(model\_name)  
accuracies.append(accuracy\_rf)  
f1s.append(f1\_rf)  
recalls.append(recall\_rf)  
recalls\_detractors.append(recall\_detractors\_rf)  
cv\_accs.append(cv\_rf)

In [290]: # Importing the adaboost packages  
from sklearn.ensemble import GradientBoostingClassifier

### 5.e) TfidfVectorizer and Gradient Boosting

# Defining the pipeline with the fixed tfidf parameters and GradientBoostingClassifier  
gb\_pipeline = imblearn.pipeline.Pipeline([  
 ('tfidf', TfidfVectorizer(stop\_words=new\_stopwords)),  
 ('model', GradientBoostingClassifier(random\_state=42)),  
 ('classifier', GradientBoostingClassifier())  
])  
Since RandomForest, an ensemble of Decision Trees, performed better than Decision Trees, I will try to take these models to the next levels with Gradient Boosting.  
# Fitting the pipeline on training data  
gb\_pipeline.fit(X\_train['lemmatized\_review'], y\_train)

Gradient boosting is a class of ensemble algorithms constructed from decision tree models. Trees

```
In [289]: # Appending each metric to the lists
model_names.append(model_name)
accuracy_rf = accuracy_rf.append(accuracy_rf)
f1s.append(f1_rf)
recalls.append(recall_rf)
recalls_detractors.append(recall_detractors_rf)
cv_acccs.append(cv_rf)

In [290]: # Importing the relevant packages
from sklearn.ensemble import GradientBoostingClassifier
```

### 5. e) TfIdfVectorizer and Gradient Boosting

```
# Defining the pipeline with the fixed tfidf parameters and GradientBoostingClass
gb_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', TfidfVectorizer(stop_words=new_stopwords)),
    ('classifier', GradientBoostingClassifier(random_state=42)),
])

1) Since RandomForest, an ensemble of Decision Trees, performed better than Decision Trees, I will
try to take these models to the next levels with Gradient Boosting.

# Fitting the pipeline on training data
gb_pipeline.fit(X_train['lemmatized_review'], y_train)

# Making predictions on test data
gb_y_pred = gb_pipeline.predict(X_test['lemmatized_review'])
```

#### 2) Evaluation Metrics

```
In [291]: # Naming the model and calling the function to evaluate it
gb_model_name = 'GradientBoosting'

# Calling the function and recording into the defined values
accuracy_gb, f1_gb, recall_gb, cv_gb = evaluation_metrics(
    y_test,
    gb_y_pred,
    gb_pipeline,
    X_train['lemmatized_review'],
    y_train
)
```

Accuracy: 0.8251  
F1-Score: 0.8313  
Recall: 0.8251  
Mean Cross-Validated Accuracy: 0.8373

```
In [292]: # Comparing the evaluation metrics
metrics_comparison('Accuracy', accuracy_rf, accuracy_gb)
metrics_comparison('F1', f1_rf, f1_gb)
metrics_comparison('Recall', recall_rf, recall_gb)
metrics_comparison('Cross-validated Accuracy', cv_rf, cv_gb)
```

Accuracy decreased from 0.8366 to 0.8251.  
F1 decreased from 0.8421 to 0.8313.

```
In [294]: # Recalling the definition of recall
gb_pipeline.fit(X_train['lemmatized_review'], y_train)
gb_y_pred = gb_pipeline.predict(X_test['lemmatized_review'])

Recall increased from 0.8366 to 0.8251.
```

```
In [293]: # Printing the overall summary
print('Random Forest actually produced better results than Gradient Boosting.')
print(classification_report(y_test, gb_y_pred))

Classification Report:
```

	precision	recall	f1-score	support
Not Detractors	0.92471	0.83011	0.87486	3773
Detractors	0.63076	0.81111	0.70966	1350
accuracy			0.82510	5123
macro avg	0.82061	0.79226	0.81233	5123
weighted avg	0.84725	0.82510	0.83133	5123

#### 4) Confusion Matrix

```
In [295]: # Recording and displaying the confusion matrix
gb_confusion_matrix = confusion_matrix_display(gb_pipeline, y_test, gb_y_pred)
```

```
F1 decreased from 0.8421 to 0.8313.
In [294]: Recall decreased from 0.8366 to 0.8251.
gb_pipeline_tuned_Accuracydecreased_from_0.8507to0.8378(y_test, gb_y_pred)
```

```
In [293]: Classification Report:
# Printing the overall summary
precision    recall   f1-score   support
print('Random Forest actually produced better results than Gradient Boosting.')

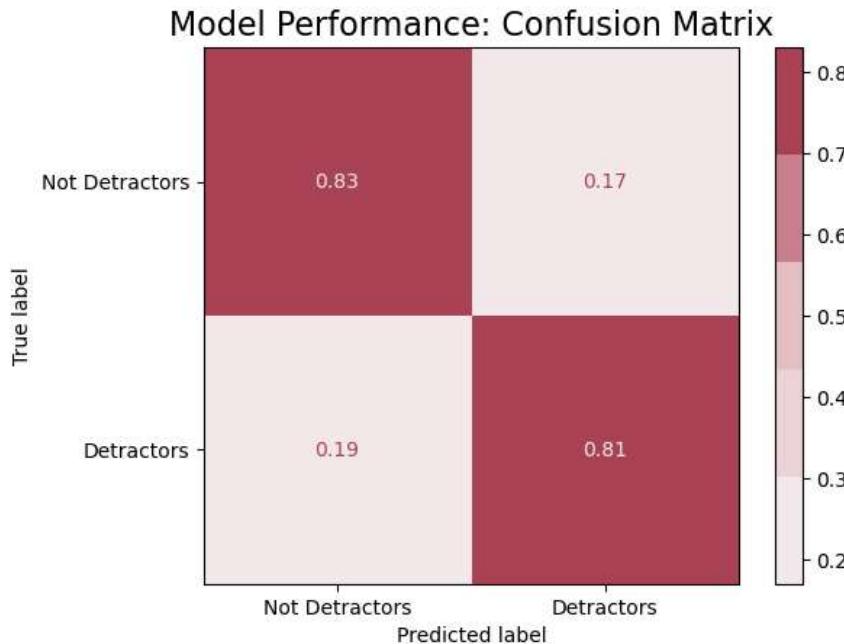
```

	Not Detractors	Detractors	
Random Forest	0.92471	0.83011	0.87486
Detractors	0.63076	0.81111	0.70966
	3773	1350	

	accuracy	precision	recall	f1-score	support
3) Classification Report	0.82061	0.82061	0.79226	0.82061	5123
weighted avg	0.84725	0.82510	0.83133	0.83133	5123

#### 4) Confusion Matrix

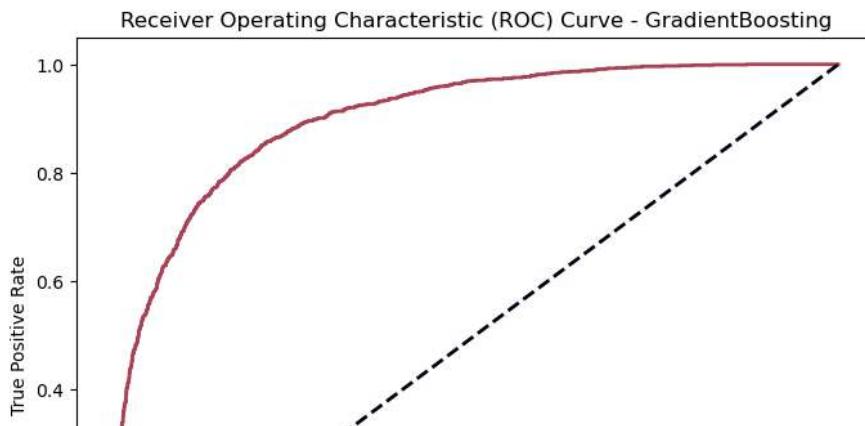
```
# Recording and displaying the confusion matrix
gb_confusion_matrix = confusion_matrix_display(gb_pipeline, y_test, gb_y_pred)
```



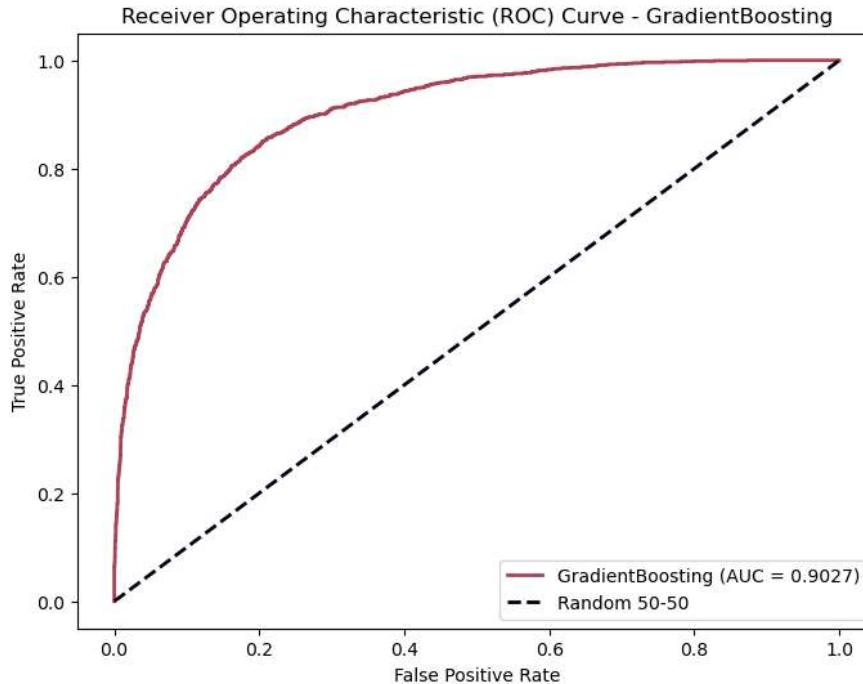
Even the detractors class got a lower true positive rate with Gradient Boosting than Random Forest.

#### 5) ROC Curve and AUC

```
# Plotting the ROC curve for the baseline model
gb_y_pred_proba = gb_pipeline.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_gb = plot_roc_curve(y_test, gb_y_pred_proba, model_name=gb_model_name)
```



```
In [296]: # Plotting the ROC curve for the baseline model
gb_y_pred_proba = gb_pipeline.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_gb = plot_roc_curve(y_test, gb_y_pred_proba, model_name=gb_model_name)
```



```
In [297]: # Printing the interpretation
auc_interpretation(auc_gb)
```

AUC = 0.9027. The model has a great predictive power in distinguishing between the positive and negative classes.

```
In [298]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_rf, auc_gb)
print()
```

Compared to the last run model:  
AUC decreased from 0.9106 to 0.9027.

```
In [299]: # Appending each metric to the lists
modelnames.append(gb_model_name)
accuracies.append(accuracy_gb)
f1s.append(f1_gb)
recalls.append(recall_gb)
recalls_detractors.append(recall_detractors_gb)
cv_accs.append(cv_gb)
roc_aucs.append(auc_gb)
```

1) Fitting and training on train data

Since Gradient Boosting did not provide better results than Random Forest, I will try another ensemble classifier.

```
In [300]: # Importing the relevant packages
from sklearn.ensemble import AdaBoostClassifier
```

## 5. f) TfidfVectorizer and AdaBoost

```
# Defining the pipeline with the fixed tfidf parameters and GradientBoostingClassifier
ab_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', TfidfVectorizer(stop_words='new_stopwords')),
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
    ('classifier', AdaBoostClassifier())
])
# While Gradient Boosting is often considered more flexible and powerful, these iterations have
# proven to provide surprising results. AdaBoost focuses on improving the classification of the
# instances that are more challenging for the current ensemble.
ab_pipeline.fit(X_train['lemmatized_review'], y_train)

# Making predictions on test data
ab_y_pred = ab_pipeline.predict(X_test['lemmatized_review'])
```

```
roc_aucs.append(auc_gb)
```

### 1) Fitting and training on train data

Since Gradient Boosting did not provide better results than Random Forest, I will try another ensemble classifier.

In [300]: # Importing the relevant packages

```
from sklearn.ensemble import AdaBoostClassifier
```

### 5. f) TfidfVectorizer and AdaBoost

```
# Defining the pipeline with the fixed tfidf parameters and GradientBoostingClass
```

```
ab_pipeline = imblearn.pipeline.Pipeline([
    ('tfidf', TfidfVectorizer(stop_words='new_stopwords')),
```

```
    ('us', RandomUnderSampler(sampling_strategy=1, random_state=42)),
```

```
    ('classifier', AdaBoostClassifier())])
```

While Gradient Boosting is often considered more flexible and powerful, these iterations have proven to provide surprising results. AdaBoost focuses on improving the classification of the instances that are more challenging for the current ensemble.

```
# Fitting the pipeline on the training data
```

```
ab_pipeline.fit(X_train['lemmatized_review'], y_train)
```

```
# Making predictions on test data
```

```
ab_y_pred = ab_pipeline.predict(X_test['lemmatized_review'])
```

### 2) Evaluation Metrics

In [301]:

```
# Naming the model and calling the function to evaluate it
ab_model_name = 'AdaBoost'
```

```
# Calling the function and recording into the defined values
```

```
accuracy_ab, f1_ab, recall_ab, cv_ab = evaluation_metrics(
    y_test,
    ab_y_pred,
    ab_pipeline,
    X_train['lemmatized_review'],
    y_train
)
```

Accuracy: 0.8150

F1-Score: 0.8221

Recall: 0.8150

Mean Cross-Validated Accuracy: 0.8276

In [302]:

```
# Comparing the evaluation metrics between undersampled model (best model) and no
metrics_comparison('Accuracy', accuracy_gb, accuracy_ab)
metrics_comparison('F1', f1_gb, f1_ab)
metrics_comparison('Recall', recall_gb, recall_ab)
metrics_comparison('Cross-validated Accuracy', cv_gb, cv_ab)
```

Accuracy decreased from 0.8251 to 0.8150.

F1 decreased from 0.8313 to 0.8221.

Recall decreased from 0.8251 to 0.8150.

Cross-validated Accuracy decreased from 0.8373 to 0.8276.

In [304]:

```
# Putting together a report
ab_pipeline_report, ab_detailed_detectors_table = classification_report(y_test, ab_y_pred)
```

AdaBoost did not provide better results than Gradient Boosting.

	precision	recall	f1-score	support
Not Detractors	0.92278	0.81712	0.86674	3773
Detractors	0.61279	0.80889	0.69732	1350
accuracy			0.81495	5123
macro avg	0.76779	0.81301	0.78203	5123
weighted avg	0.84109	0.81495	0.82210	5123

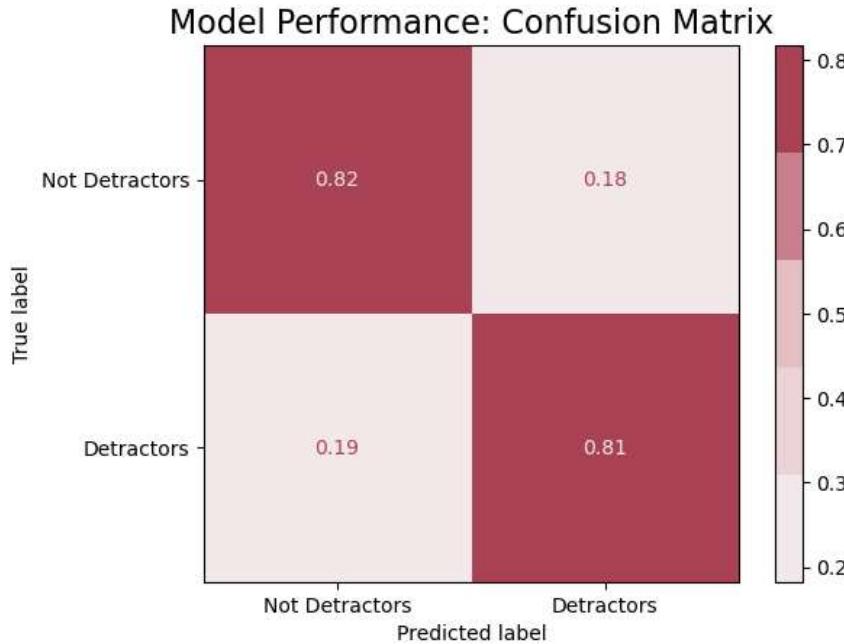
```
In [304]: # Path to the report summary
ab_pipeline_report, ab_recall_detractors_tab = classsification(y_test, ab_y_pred)
```

AdaBoost classifier provides better results than Gradient Boosting.

	precision	recall	f1-score	support
Not Detractors	0.92278	0.81712	0.86674	3773
Detractors	0.61279	0.80889	0.69732	1350
accuracy			0.81495	5123
macro avg	0.76779	0.81301	0.78203	5123
weighted avg	0.84109	0.81495	0.82210	5123

#### 4) Confusion Matrix

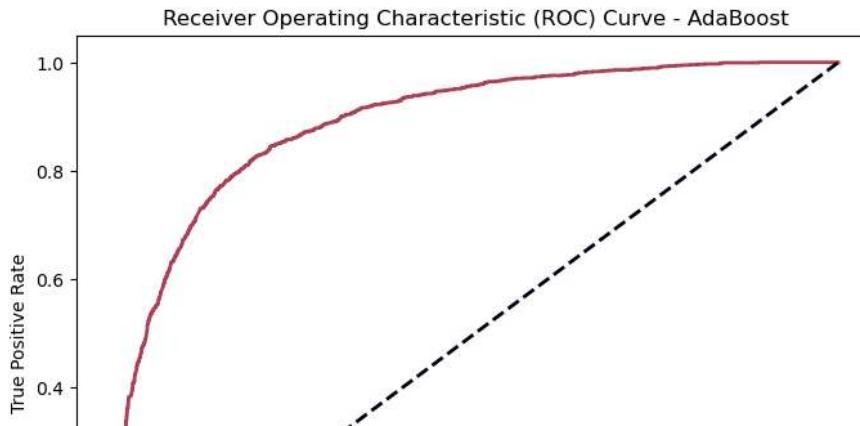
```
In [305]: # Recording and displaying the confusion matrix
ab_confusion_matrix = confusion_matrix_display(ab_pipeline, y_test, ab_y_pred)
```



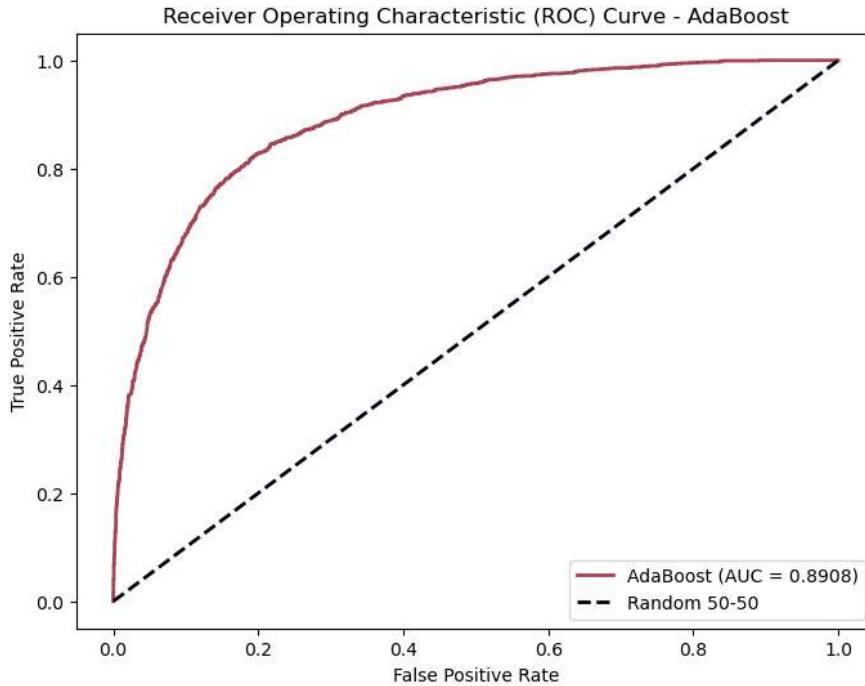
The lower results are also reflected on recall for detractors.

#### 5) ROC Curve and AUC

```
In [306]: # Plotting the ROC curve for the baseline model
ab_y_pred_proba = ab_pipeline.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_ab = plot_roc_curve(y_test, ab_y_pred_proba, model_name=ab_model_name)
```



```
In [306]: # Plotting the ROC curve for the baseline model
ab_y_pred_proba = ab_pipeline.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_ab = plot_roc_curve(y_test, ab_y_pred_proba, model_name=ab_model_name)
```



```
In [307]: # Printing the interpretation
auc_interpretation(auc_ab)
```

AUC = 0.8908. The model has a good predictive power.

```
In [308]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_gb, auc_ab)
print()

# # Interpreting AUC:
# print('Compared to the best model so far: ')
# metrics_comparison('AUC', auc_resamp1, auc_prep)
```

Compared to the last run model:  
AUC decreased from 0.9027 to 0.8908.

```
In [309]: # Appending each metric to the lists
modelnames.append(ab_model_name)
accuracies.append(accuracy_ab)
f1s.append(f1_ab)
recalls.append(recall_ab)
recalls_detractors.append(recall_detractors_ab)
cv_accs.append(cv_ab)
roc_aucs.append(auc_ab)
```

```
In [310]: # Importing the relevant packages
from xgboost import XGBClassifier
```

### 5. g) TfidfVectorizer and XGBoost

```
# Defining the pipeline with the fixed tfidf parameters and GradientBoostingClassifier
15th iteration XGBoost TfidfVectorizer
('tfidf', TfidfVectorizer(stop_words=new_stopwords)),
('us', RandomUnderSampler(sampling_strategy=1, random_state=42))
Since I am not reaching stronger results than with NaiveBayes with the previous classifiers, I will
try XGBoost, for its popularity among classifiers algorithms.
```

XGBoost, or Extreme Gradient Boosting, is a powerful and efficient machine learning algorithm that belongs to the gradient boosting family, known for its speed, scalability, and ability to handle diverse data types.

```
# Making predictions on test data
xgb_y_pred = xgb_pipeline.predict(X_test['lemmatized_review'])
```

```

recalls.append(recall_ab)
recalls_detractors.append(recall_detractors_ab)
    1) Fitting and training on train data
cv_accs.append(cv_ab)
roc_aucs.append(auc_ab)

```

In [310]: # Importing the relevant packages

### 5. g) TfidfVectorizer and XGBoost

# Defining the pipeline with the fixed tfidf parameters and GradientBoostingClassifier

15th iteration XGBoost, Tfidf Vectorizer

('tfidf', TfidfVectorizer(stop\_words=new\_stopwords)),

('us', RandomUnderSampler(sampling\_strategy=1, random\_state=42)),

Since I am not reaching stronger results than with NaiveBayes with the previous classifiers, I will try XGBoost, for its popularity among classifiers algorithms.

# Fitting the pipeline on training data  
that helps to the gradient boosting family known for its speed, scalability, and ability to handle diverse data types.

# Making predictions on test data

xgb\_y\_pred = xgb\_pipeline.predict(X\_test['lemmatized\_review'])

## 2) Evaluation Metrics

In [311]: # Naming the model and calling the function to evaluate it  
xgb\_model\_name = 'XGBoost'

# Calling the function and recording into the defined values

```

accuracy_xgb, f1_xgb, recall_xgb, cv_xgb = evaluation_metrics(
    y_test,
    xgb_y_pred,
    xgb_pipeline,
    X_train['lemmatized_review'],
    y_train
)

```

Accuracy: 0.8468

F1-Score: 0.8520

Recall: 0.8468

Mean Cross-Validated Accuracy: 0.8574

In [312]: # Comparing the evaluation metrics between the last model the current one

```

metrics_comparison('Accuracy', accuracy_ab, accuracy_xgb)
metrics_comparison('F1', f1_ab, f1_xgb)
metrics_comparison('Recall', recall_ab, recall_xgb)
metrics_comparison('Cross-validated Accuracy', cv_ab, cv_xgb)

```

Accuracy improved from 0.8150 to 0.8468.

F1 improved from 0.8221 to 0.8520.

Recall improved from 0.8150 to 0.8468.

Cross-validated Accuracy improved from 0.8276 to 0.8574.

In [314]: # Putting the final report summary

xgb\_pipeline\_report = classification\_report(y\_test, xgb\_y\_pred)

Classification Report can be confirmed! Results highly improved.  
precision      recall      f1-score      support

	Not Detractors	Detractors	accuracy	macro avg	weighted avg
3) Classification Report	0.93812	0.84787	0.84677	0.80152	0.86613
	0.66492	0.84370	0.74372	0.84579	0.84677
			3773	5123	5123
			1350		

## 4) Confusion Matrix

In [315]: # Recording and displaying the confusion matrix

xgb\_confusion\_matrix = confusion\_matrix\_display(xgb\_pipeline, y\_test, xgb\_y\_pred)

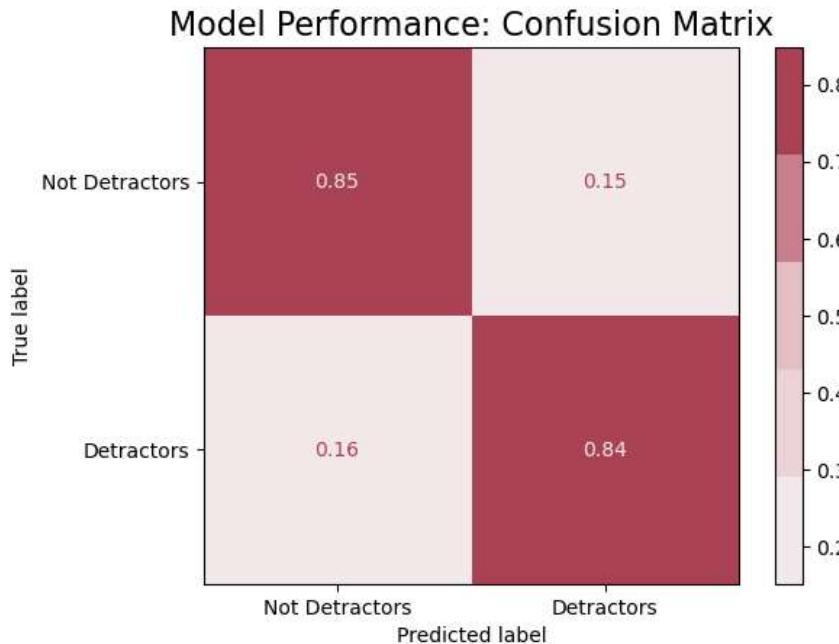
```
In [314]: # Path to the previous summary
#gb_dt_clf_xgb_report_topically_detractors_xgbmed.classification(y_test, xgb_y_pred)

# Baseline report can be confirmed! Results highly improved.
precision    recall   f1-score   support
Not Detractors      0.93812   0.84787   0.89071   3773
Detractors          0.66492   0.84370   0.74372   1350

accuracy           0.84677
macro avg       0.80152   0.84579   0.81721   5123
weighted avg     0.86613   0.84677   0.85198   5123
```

## 4) Confusion Matrix

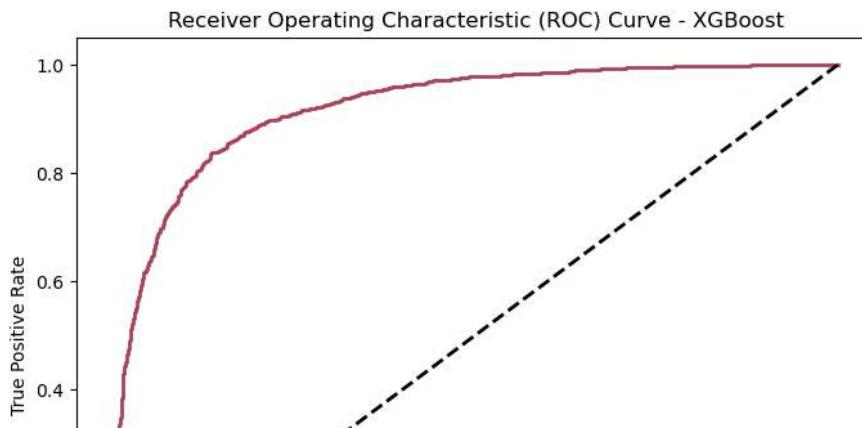
```
In [315]: # Recording and displaying the confusion matrix
xgb_confusion_matrix = confusion_matrix_display(xgb_pipeline, y_test, xgb_y_pred)
```



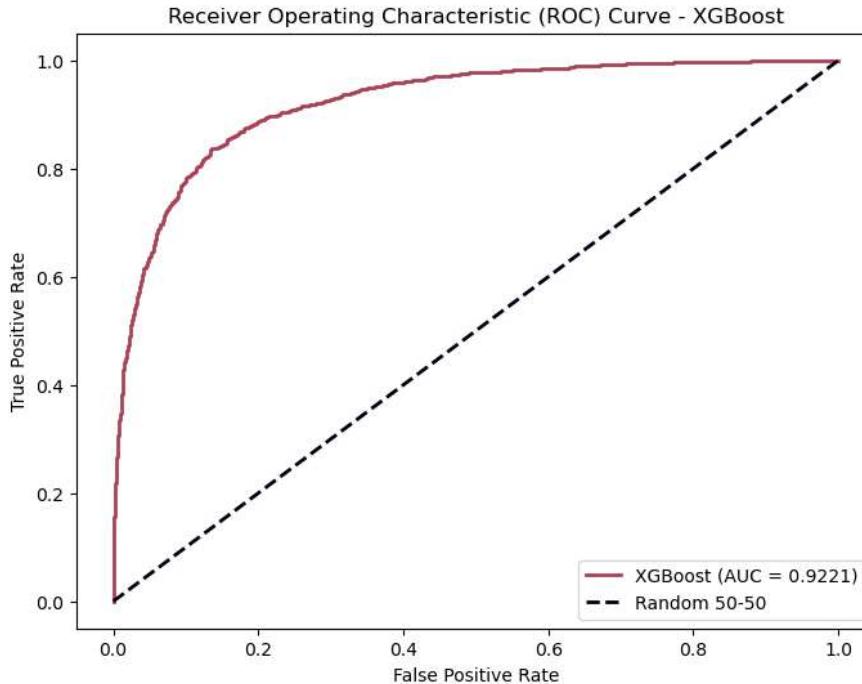
Nevertheless, recall for detractors is still not reaching higher results than with the Naive Bayes model.

## 5) ROC Curve and AUC

```
In [316]: # Plotting the ROC curve for the baseline model
xgb_y_pred_proba = xgb_pipeline.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_xgb = plot_roc_curve(y_test, xgb_y_pred_proba, model_name=xgb_model_name)
```



```
In [316]: # Plotting the ROC curve for the baseline model
xgb_y_pred_proba = xgb_pipeline.predict_proba(X_test['lemmatized_review'])[:, 1]
auc_xgb = plot_roc_curve(y_test, xgb_y_pred_proba, model_name=xgb_model_name)
```



```
In [317]: # Printing the interpretation
auc_interpretation(auc_xgb)
```

AUC = 0.9221. The model has a great predictive power in distinguishing between the positive and negative classes.

```
In [318]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_ab, auc_xgb)
print()
```

Compared to the last run model:  
AUC improved from 0.8908 to 0.9221.

```
In [319]: # Appending each metric to the lists
modelnames.append(xgb_model_name)
accuracies.append(accuracy_xgb)
f1s.append(f1_xgb)
recalls.append(recall_xgb)
recalls_detractors.append(recall_detractors_xgb)
cv_accs.append(cv_xgb)
```

```
In [320]: # Defining the parameter grid
param_grid = {
    'estimator__max_depth': [3, 5, 7],
    'estimator__subsample': [0.8, 0.9, 1.0],
}
# 16th iteration: Tuned XGBoost
```

I finally got another classifier competing with the results provided from Multinomial Naive Bayes. I will try to tune it and find the best parameters to see if these results can be beaten.

```
In [321]: # Creating the GridSearchCV object
grid_search = GridSearchCV(xgb_pipeline, param_grid, cv=5, scoring='recall', n_jobs=-1)

# Fitting the grid search to the data
grid_search.fit(X_train['lemmatized_review'], y_train)
```

With try to apply combinatorial grid searching for XGBoost to get the best parameters and try to improve further the model.

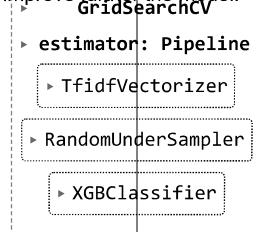
```
Out[321]: GridSearchCV
  estimator: Pipeline
    > TfidfVectorizer
      > RandomizedSampler
```

```
In [320]: # Defining the parameter grid
param_grid = {
    'classifier_n_estimators': [50, 100, 150],
    'classifier_max_depth': [3, 5, 7],
    'classifier_subsample': [0.8, 0.9, 1.0],
}
16th iteration: Tuned XGBoost
```

```
In [321]: # Creating the GridSearchCV object
grid_search = GridSearchCV(xgb_pipeline, param_grid, cv=5, scoring='recall', n_jobs=-1)

# Fitting the grid search to the data
grid_search.fit(X_train['lemmatized_review'], y_train)
```

Out[321]:



```
In [322]: # Storing the best params
best_xgb_params = grid_search.best_params_

# Storing the best model from the grid search
best_xgb_pipeline = grid_search.best_estimator_

# Printing the best parameters and the corresponding score
print("Best Parameters: ", grid_search.best_params_)
```

Best Parameters: {'classifier\_learning\_rate': 0.2, 'classifier\_max\_depth': 7, 'classifier\_n\_estimators': 150}

```
In [323]: # Fitting the best pipeline on training data
best_xgb_pipeline.fit(X_train['lemmatized_review'], y_train)

# Making predictions on the test set
best_xgb_y_pred = best_xgb_pipeline.predict(X_test['lemmatized_review'])
```

## 2) Evaluation Metrics

```
In [324]: # Naming the model and calling the function to evaluate it
xgbtuned_model_name = 'Tuned XGBoost'
```

```
# Calling the function and recording into the defined values
```

```
In [325]: # Comparing the evaluation metrics between the last model the current one
y_test = X_test['lemmatized_review']
accuracy_xgbtuned, f1_xgbtuned, recall_xgbtuned, cv_xgbtuned = evaluation_metrics(
    y_test,
    metrics_comparison('Accuracy', accuracy_xgb, accuracy_xgbtuned),
    metrics_comparison('F1', f1_xgb, f1_xgbtuned),
    metrics_comparison('Recall', recall_xgb, recall_xgbtuned),
    X_train['lemmatized_review'],
    metrics_comparison('Cross-Validated Accuracy', cv_xgb, cv_xgbtuned),
    y_train
)
```

Accuracy improved from 0.8468 to 0.8518.  
 Accuracy: 0.8518  
 F1 improved from 0.8520 to 0.8565.  
 F1-Score: 0.8565  
 Recall improved from 0.8468 to 0.8518.  
 Recall: 0.8518  
 Cross-Validated Accuracy decreased from 0.8574 to 0.8567.  
 Mean Cross-Validated Accuracy: 0.8567

```
In [326]: # Printing the overall summary
print('Despite the expensive computing, overall scores are the best ones recorded decreased compared to the untuned XGBoost model but remains over 85%.)
```

Despite the expensive computing, overall scores are the best ones recorded. Cross-validated accuracy slightly decreased compared to the untuned XGBoost model but remains over 85%.

```
In [325]: # Comparing the evaluation metrics between the last model the current one
accuracy_xgbtuned, f1_xgbtuned, recall_xgbtuned, cv_xgbtuned = evaluation_metrics(y_test,
best_xgb_y_pred, metrics_comparison('Accuracy', accuracy_xgb, accuracy_xgbtuned)
metrics_comparison('F1', f1_xgb, f1_xgbtuned)
metrics_comparison('Recall', recall_xgb, recall_xgbtuned)
metrics_comparison('Cross-validated review', cv_xgb, cv_xgbtuned)
y_train
y_train)
```

Accuracy improved from 0.8468 to 0.8518.  
Accuracy: 0.8518  
F1 Score: 0.8565  
Recall improved from 0.8468 to 0.8518.  
Recall: 0.8518  
Cross-Validated Accuracy decreased from 0.8574 to 0.8567.  
Mean Cross-Validated Accuracy: 0.8567

```
In [326]: # Printing the overall summary
print('Despite the expensive computing, overall scores are the best ones recorded decreased compared to the untuned XGBoost model but remains over 85%.' )
```

Despite the expensive computing, overall scores are the best ones recorded. Cross-validated accuracy slightly decreased compared to the untuned XGBoost model but remains over 85%.

### 3) Classification Report

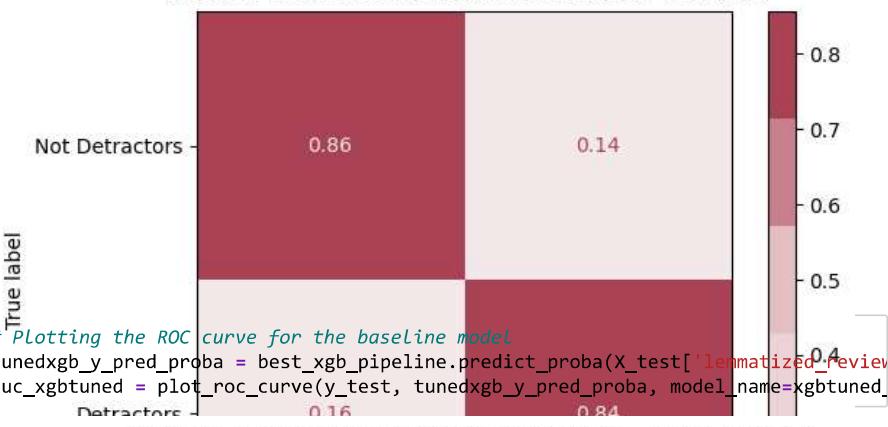
```
In [327]: # Calling the function
tunedxgb_class_report, recall_detractors_tunedxgb = class_calculation(y_test, best_xgb_y_pred)
```

Classification Report:				
	precision	recall	f1-score	support
Not Detractors	0.93706	0.85635	0.89489	3773
Detractors	0.67642	0.83926	0.74909	1350
accuracy			0.85184	5123
macro avg	0.80674	0.84780	0.82199	5123
weighted avg	0.86838	0.85184	0.85647	5123

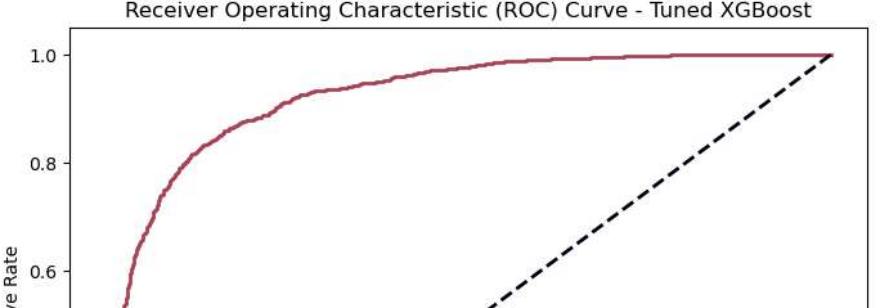
### 4) Confusion Matrix

```
In [328]: # Recording and displaying the confusion matrix
tunedxgb_confusion_matrix = confusion_matrix_display(best_xgb_pipeline, y_test, tunedxgb_y_pred)
```

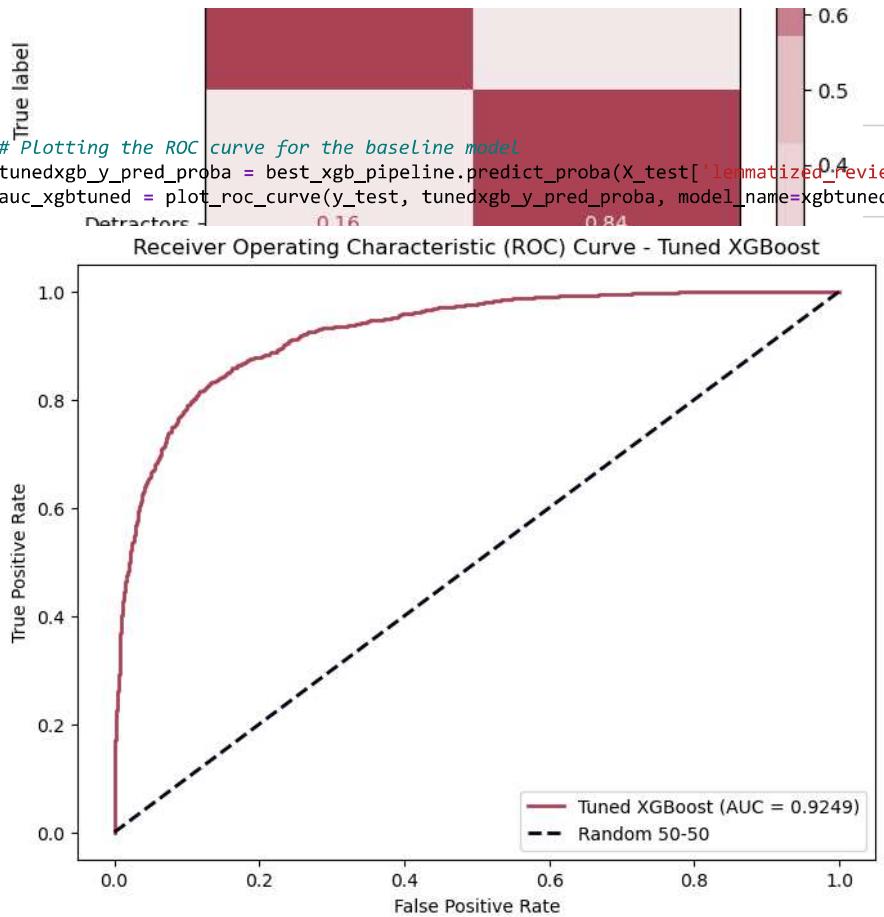
Model Performance: Confusion Matrix



```
In [329]: # Plotting the ROC curve for the baseline model
tunedxgb_y_pred_proba = best_xgb_pipeline.predict_proba(X_test['lemmatized_review'])
auc_xgbtuned = plot_roc_curve(y_test, tunedxgb_y_pred_proba, model_name=xgbtuned)
```



```
In [329]: # Plotting the ROC curve for the baseline model
tunedxgb_y_pred_proba = best_xgb_pipeline.predict_proba(X_test['lemmatized_review'])
auc_xgbtuned = plot_roc_curve(y_test, tunedxgb_y_pred_proba, model_name=xgbtuned)
```



```
In [330]: # Printing the interpretation
auc_interpretation(auc_xgbtuned)
```

AUC = 0.9249. The model has a great predictive power in distinguishing between the positive and negative classes.

```
In [331]: # Interpreting AUC:
print('Compared to the last run model: ')
metrics_comparison('AUC', auc_xgb, auc_xgbtuned)
print()
```

Compared to the last run model:  
AUC improved from 0.9221 to 0.9249.

```
In [332]: # Appending each metric to the lists
modelnames.append(xgbtuned_model_name)
accuracies.append(accuracy_xgbtuned)
f1s.append(f1_xgbtuned)
cv_acccs.append(cv_xgbtuned)
recalls_detractors.append(recall_detractors_tunedxgb)
cv_accs.append(cv_xgbtuned)
roc_aucs.append(auc_xgbtuned)
```

## 6. Evaluation

### 6. a) Final Model and Classification Metrics

```
In [333]: # Verifying that all lists have the same size
# print(len(modelnames))
# print()
# print(len(accuracies))
# print()
# print(len(f1s))
# print()
# print(len(recalls))
# print()
# print(len(recalls_detractors))
# print()
# print(len(cv_acccs))
# print()
# print(len(roc_aucs))
```

```

    accuracies.append(accuracy_xgbtuned)
    f1s.append(f1_xgbtuned)
    recall_detractors.append(recall_detractors_tunedxgb)
    cv_accs.append(cv_xgbtuned)
    roc_aucs.append(auc_xgbtuned)
6. a) Final Model and Classification Metrics

```

In [333]:

```

# Verifying that all lists have the same size
# print(len(modelnames))
# print()
# print(len(accuracies))
# print()
# print(len(f1s))
# print()
# print(len(recalls))
# print()
# print(len(recalls_detractors))
# print()
# print(len(cv_accs))
# print()
# print(len(roc_aucs))
# print()

```

In [334]:

```

# Creating a DataFrame with stored best scores

models_results = {
    'models': modelnames,
    'accuracy': accuracies,
    'f1': f1s,
    'recall': recalls,
    'recall_detractors': recalls_detractors,
    'cross-val_accuracies': cv_accs,
    'ROC-AUCs': roc_aucs
}

overall = pd.DataFrame(models_results)

```

In [335]:

```

# Inspecting all iterations' results
overall

```

Out[335]:

	models	accuracy	f1	recall	recall_detractors	cross-val_accuracies	ROC-AUCs
0	Baseline	0.754636	0.665579	0.754636	0.06963	0.750716	0.911003
1	Resampled	0.884248	0.878426	0.884248	0.65481	0.883589	0.927415
2	Sampling 1	0.831154	0.837767	0.831154	0.84296	0.844221	0.919755
3	No Stopwords	0.831934	0.838196	0.831934	0.83333	0.846628	0.918632
4	Industry Stopwords	0.833106	0.839580	0.833106	0.84444	0.845977	0.920165
5	Full Stopwords	0.834472	0.840474	0.834472	0.83259	0.848060	0.919093
6	Lemmatized	0.827054	0.834410	0.827054	0.85630	0.838169	0.919214
7	Preprocessed	0.827054	0.834132	0.827054	0.84593	0.840381	0.917364

In [337]:

```

# Printing metrics for max recall index
overall.iloc[max_recall_detractors_index]

```

Out[337]:

8	NaiveBayes GS	0.827445	0.834899	0.827445	0.86148	0.836998	0.914550
9	DecisionTree	0.729260	0.742067	0.729260	0.70519	0.724426	0.721530
10	RandomForest	0.836619	0.842079	0.836619	0.82148	0.850663	0.910567
11	GradientBoosting	0.825102	0.801327	0.825102	0.81111	0.837323	0.902704
12	AdaBoost	0.814952	0.822098	0.814952	0.80889	0.827563	0.890777
13	cross-val_accuracies		0.86148				
14	XGBoost	0.846769	0.850791	0.846769	0.84370	0.857366	0.922079
15	Name: 8, dtype: object	0.851845	0.856469	0.851845	0.83926	0.856715	0.924855

In [339]:

```

# Storing both recalls the columns I wanted to average
# Storing the maximum recall detractors value and its index
both_recalls = [recall_detractors]
max_recall_detractor = max(recalls_detractors)
max_recall_index_detractor = recalls_detractors.index(max_recall_detractor)
# Calculate the average for each row
overall['average_both_recalls'] = overall[both_recalls].mean(axis=1)

```

In [340]:

```
# overall
```

```
In [337]: # Printing metrics for max recall index
overall.iloc[max_recall_index].values[145071]
```

8	NaiveBayes GS	0.827445	0.834899	0.827445	0.86148	0.836998	0.914550
9	DecisionTree	0.729260	NaiveBayes GS	0.729260	0.72815	0.756442	0.803643
10	accuracy	0.827445	0.836619	0.827445	0.70519	0.724426	0.721530
11	RandomForest	0.836619	0.827445	0.836619	0.82148	0.850663	0.910567
12	f1	0.834899	0.834899	0.834899	0.81111	0.837323	0.902704
13	GradientBoosting	0.825102	0.8038674	0.825102	0.80889	0.827563	0.890777
14	recall_detractors	0.86148	0.86148	0.86148	0.84370	0.857366	0.922079
15	AdaBoost	0.814952	0.822096	0.814952	0.836998	0.856715	0.924855
Name: 8	dtype: object						
16	Tuned XGBoost	0.851845	0.856469	0.851845	0.83926	0.856715	0.924855

```
In [339]: # Storing both recalls the columns I wanted to average
In [340]: # Storing the maximum recall detractors value and its index
both_recalls = recall_detractors
max_recall_detractor = max(both_recalls)
max_recall_index_detractor = recall_detractors.index(max_recall_detractor)
# Calculate the average for each row
overall['average_both_recalls'] = overall[both_recalls].mean(axis=1)
```

```
In [340]: # overall
```

```
In [341]: # Storing the maximum recall for both recalls and its index
max_both_recalls = max(both_recalls)
max_both_recalls_index = both_recalls.index(max_both_recalls)
```

```
In [342]: # Creating a bar chart to review all
```

```
# Creating a bar chart for recall
fig, ax = plt.subplots(nrows=1, ncols=1, figsize=(10,6))

# Setting facecolor
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')

# plt.figure(figsize=(10, 6))
plt.bar(overall['models'], overall['recall_detractors'], color='#AA4255', label='Recall Detractors')

# Creating Line plots for Log Loss and accuracy
plt.plot(overall['models'], overall['accuracy'], marker='o', color="#00AF87", label='Accuracy')
# plt.plot(overall['models'], overall['f1'], marker='x', color="#00A4EF", label='F1 Score')
plt.plot(overall['models'], overall['ROC-AUCs'], marker='x', color="#00A4EF", label='ROC-AUCs')

# Setting Labels and title
plt.xlabel('Models', fontsize=12, color='white')
plt.xticks(rotation=30, color='white')
...
```

```
In [342]: # Creating a bar chart to review all

# Creating a bar chart for recall
fig, ax = plt.subplots(nrows=1, ncols=1, figsize=(10,6))

# Setting facecolor
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')

# plt.figure(figsize=(10, 6))
plt.bar(overall['models'], overall['recall_detractors'], color='#AA4255', label='Recall')

# Creating line plots for Log Loss and accuracy
plt.plot(overall['models'], overall['accuracy'], marker='o', color='#00AF87', label='Accuracy')
# plt.plot(overall['models'], overall['f1'], marker='x', color='#00A4EF', label='F1 Score')
plt.plot(overall['models'], overall['ROC-AUCs'], marker='x', color='#00A4EF', label='ROC-AUCs')

# Setting Labels and title
plt.xlabel('Models', fontsize=12, color='white')
plt.xticks(rotation=30, color='white')
ax.tick_params(axis='y', colors='white')

plt.ylabel('recall, accuracy, CV, AUC', fontsize=12, color='white')

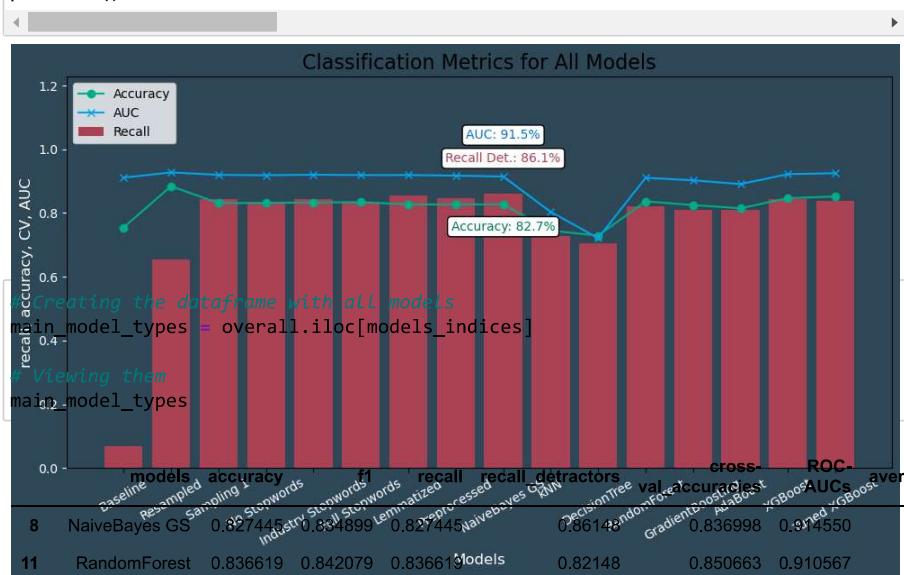
# Annotating the last index of each category
plt.annotate(f'Recall Det.: {recalls_detractors[max_recall_index_detrac] * 100:.1f}%', (modelnames[-1], recalls_detractors[max_recall_index_detrac] * 100))
plt.annotate(f'Accuracy: {accuracies[max_recall_index_detrac] * 100:.1f}%', (modelnames[-1], accuracies[max_recall_index_detrac] * 100))
plt.annotate(f'AUC: {roc_auc[0]:.1f}%', (modelnames[-1], roc_auc[0]))

# Defining the max value of y
max_y = max(overall[['ROC-AUCs', 'recall_detractors']].max())
plt.ylim(0, max_y + 0.3)
plt.title('Classification Metrics for All Models', fontsize=16)

# Displaying the legend
plt.legend(loc='upper left')

# Saving the plot as a PNG with a transparent background
plt.savefig('images/final_model_names.png', transparent=True)

# Showing the plot
plt.tight_layout()
plt.show()
```



```
In [344]: #Creating the dataframe with all models
main_model_types = overall.iloc[models_indices]

#Viewing them
main_model_types
```

```
Out[344]:
```

Index	Model Type	NaiveBayes GS	Resampled accuracy	Sampling Stewards	Industry Stewards	Lemmatized Stewards	stopword-processed	NaiveBayes ANN	DecisionTree	randomFor	val	cross-validated	ROC-AUCs	average
8	NaiveBayes GS	0.827445	0.836619	0.834899	0.834899	0.827445	0.827445	0.86148	0.836998	0.850663	0.845455	0.84370	0.84370	
11	RandomForest	0.825102	0.831326	0.832079	0.836619	0.834899	0.834899	0.82148	0.837323	0.850663	0.910567	0.84370	0.84370	
12	GradientBoosting	0.825102	0.831326	0.832079	0.836619	0.834899	0.834899	0.81111	0.837323	0.850663	0.902704	0.84370	0.84370	
14	XGBoost	0.846769	0.851977	0.846769	0.851977	0.846769	0.846769	0.84370	0.857366	0.857366	0.922079	0.84370	0.84370	

```
In [343]: # Defining indices to filter on
models_indices = [8, 11, 12, 14]
```

```
In [344]: #Creating the dataframe with all models
main_model_types = overall.iloc[models_indices]

# Viewing them
main_model_types
```

	models	accuracy	f1	recall	recall	detractors	val	cross	ROC-AUCs	average
8	NaiveBayes GS	0.827445	0.834899	0.827445	0.827445	0.86148	0.86148	0.86148	0.836998	0.836998
11	RandomForest	0.836619	0.842079	0.836619	0.836619	0.82148	0.82148	0.82148	0.850663	0.850663
12	GradientBoosting	0.825102	0.831326	0.825102	0.825102	0.81111	0.81111	0.81111	0.837323	0.837323
14	XGBoost	0.846769	0.851977	0.846769	0.846769	0.84370	0.84370	0.84370	0.857366	0.857366

Summarizing the best classification model types

```
In [343]: # Defining indices to filter on
```

```
models_indices = [8, 11, 12, 14]
```

```
In [345]: # Creating a bar chart to review all

# Creating a bar chart for recall
fig, ax = plt.subplots(nrows=1, ncols=1, figsize=(8,6))

# Setting facecolor
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')

# plt.figure(figsize=(10, 6))
plt.bar(main_model_types['models'], main_model_types['recall_detractors'], color='red')

# Creating Line plots for Log Loss and accuracy
plt.plot(main_model_types['models'], main_model_types['accuracy'], marker='o', color='blue')
# plt.plot(main_model_types['models'], main_model_types['f1'], marker='x', color='green')
plt.plot(main_model_types['models'], main_model_types['ROC-AUCs'], marker='x', color='red')

# Setting Labels and title
plt.xlabel('Models', fontsize=12, color='white')
plt.xticks(rotation=45, color='white', ha='right')
...
```

```
In [345]: # Creating a bar chart to review all

# Creating a bar chart for recall
fig, ax = plt.subplots(nrows=1, ncols=1, figsize=(8,6))

# Setting facecolor
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')

# plt.figure(figsize=(10, 6))
plt.bar(main_model_types['models'], main_model_types['recall_detractors'], color='red')

# Creating line plots for Log Loss and accuracy
plt.plot(main_model_types['models'], main_model_types['accuracy'], marker='o', color='green')
# plt.plot(main_model_types['models'], main_model_types['f1'], marker='x', color='blue')
plt.plot(main_model_types['models'], main_model_types['ROC-AUCs'], marker='x', color='cyan')

# Setting labels and title
plt.xlabel('Models', fontsize=12, color='white')
plt.xticks(rotation=45, color='white', ha='right')
ax.tick_params(axis='y', colors='white')

plt.ylabel('Recall, accuracy, AUC', fontsize=12, color='white')

# Annotating the last index of each category
plt.annotate(f'Recall Det.: {recalls_detractors[max_recall_index_detrac] * 100:.2f}%', (modelnames[-1], recalls_detractors[-1]), color='white')
plt.annotate(f'Accuracy: {accuracies[max_recall_index_detrac] * 100:.2f}%', (modelnames[-1], accuracies[-1]), color='white')
plt.annotate(f'AUC: {roc_aucss[max_recall_index_detrac] * 100:.2f}%', (modelnames[-1], roc_aucss[-1]), color='white')

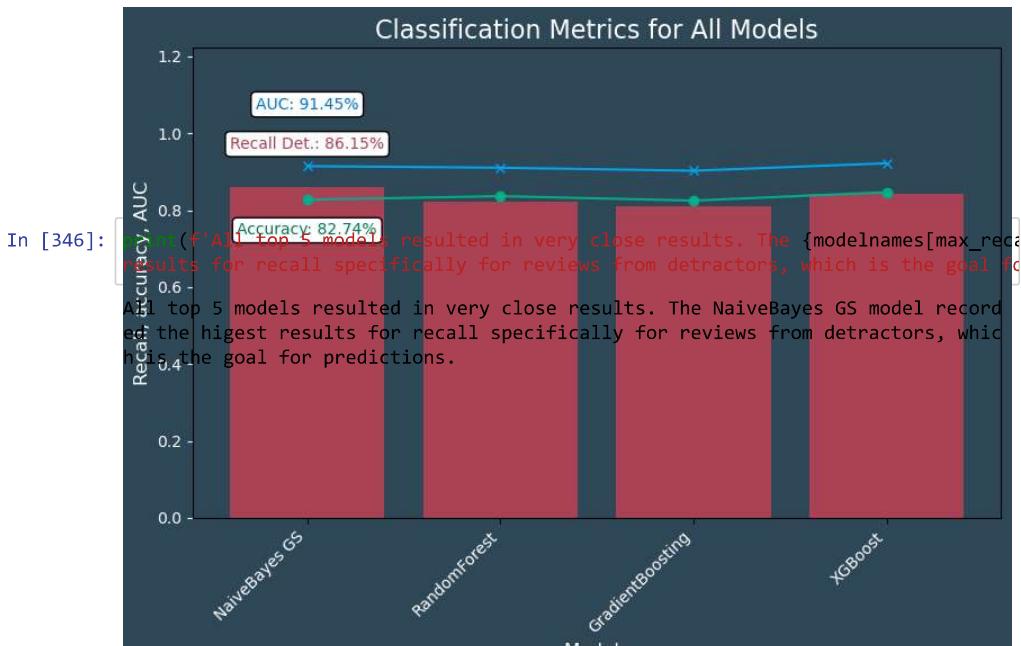
# Defining the max value of y
max_y = max(main_model_types[['ROC-AUCs', 'recall_detractors']].max())
plt.ylim(0, max_y + 0.3)
plt.title('Classification Metrics for All Models', fontsize=16, color='white')

# Displaying the legend
# plt.legend(loc='upper left')

# Showing the plot
plt.tight_layout()

# Saving the plot as a PNG
plt.savefig('images/final_models.png')

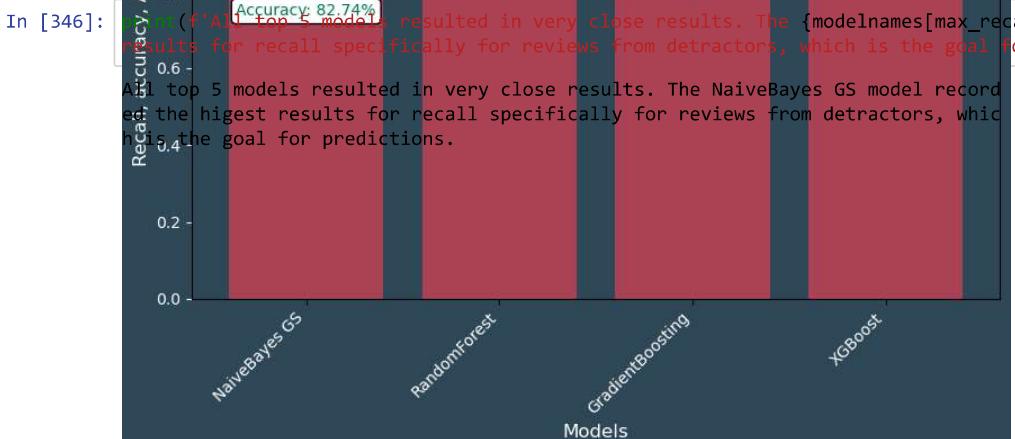
plt.show()
```



```
In [346]: print(f'All top 5 models resulted in very close results. The {modelnames[max_recall_index_detrac]} model recorded the highest results for recall specifically for reviews from detractors, which is the goal for predictions.)
```

```
All top 5 models resulted in very close results. The NaiveBayes GS model recorded the highest results for recall specifically for reviews from detractors, which is the goal for predictions.
```

```
In [347]: # Naming the model and calling the function to evaluate it
best_model_name = modelnames[max_recall_index_detrac]
best_model = best_pipeline
```



In [347]:

```
# Naming the model and calling the function to evaluate it
best_model_name = modelnames[max_recall_index_detrac]
best_model = best_pipeline

# Reminding the evaluation metrics for the model
print(f'Evaluation Metrics on Train Data for {best_model_name} model.')

print(f'Accuracy: {overall.iloc[max_recall_index_detrac]["accuracy"]:.4f}')
print(f'F1-Score: {overall.iloc[max_recall_index_detrac]["f1"]:.4f}')
print(f'Recall: {overall.iloc[max_recall_index_detrac]["recall"]:.4f}')
print(f'Recall Detractors: {overall.iloc[max_recall_index_detrac]["recall_detractors"]:.4f}')
print(f'Mean Cross-Validated Accuracy: {overall.iloc[max_recall_index_detrac]["cv_accuracy"]:.4f}'
```

Evaluation Metrics on Train Data for NaiveBayes GS model.

```
Accuracy: 0.8274
F1-Score: 0.8349
Recall: 0.8274
Recall Detractors: 0.8615
Mean Cross-Validated Accuracy: 0.8370
```

- Test data

As I have evaluated the models, I will run our function again. The same results are expected for all metrics, only the cross-validated accuracy should change as this one will be evaluated on the test, unseen data.

In [348]:

```
# Running evaluation metrics on test data
print(f'Evaluation Metrics on Unseen Data for {best_model_name} model.')
print()

# Calculating predictions using this model
best_y_pred = best_model.predict(X_test['lemmatized_review'])

# Calling the function and recording into the defined values
accuracy_best_test, recall_best_test, f1_best_test, cv_best_test = evaluation_metrics(y_test,
    best_y_pred, recall_detractors_best = class_calculation(y_test, best_y_pred,
    X_test['lemmatized_review'],
    Classification Report:
    y_test
) precision recall f1-score support
```

In [349]:

```
# Calling the function
best_y_pred, recall_detractors_best = class_calculation(y_test, best_y_pred,
    X_test['lemmatized_review'],
    Classification Report:
    y_test
) precision recall f1-score support
```

```
Model Metrics for unseen data for NaiveBayes GS model.
      Detractors  0.62527  0.86148  0.72461   1350
Accuracy: 0.8274
F1-Score: 0.8349
Recall: 0.8274
Mean Cross-Validated Accuracy: 0.8370
```

- 3) Classification Report
- 4) Confusion Matrix

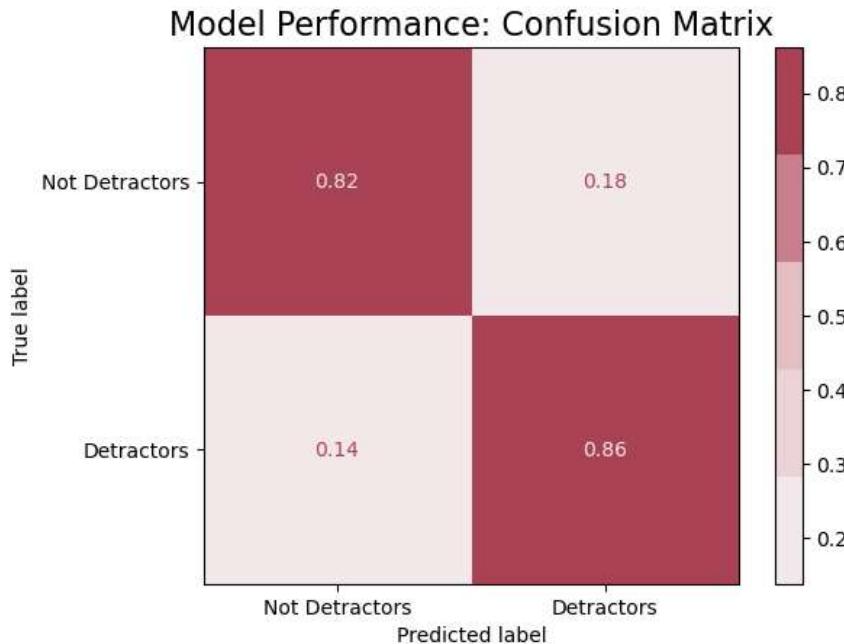
In [350]:

```
# Recording and displaying the confusion matrix
best_confusion_matrix = confusion_matrix_display(best_model, y_test, best_y_pred)
```

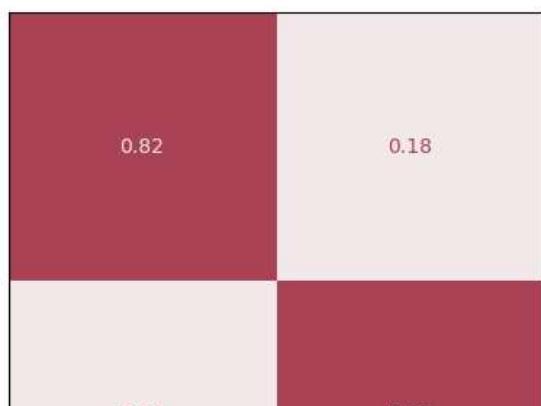
```
In [349]: # Call best_y_pred function
best_yes_model_report, recall_detractors_best = class_calculation(y_test, best_y_pred,
X_test['lemmatized_review'],
Classification Report:
y_test      precision    recall   f1-score   support
)
          precision    recall   f1-score   support
Detractors       0.62527   0.86148   0.72461    1350
Accuracy: 0.8274
F1-Score: 0.82749
Recall: 0.82749
Mac F1 Score: 0.82749
Weighted Accuracy: 0.82749
```

- 3) Classification Report  
 4) Confusion Matrix

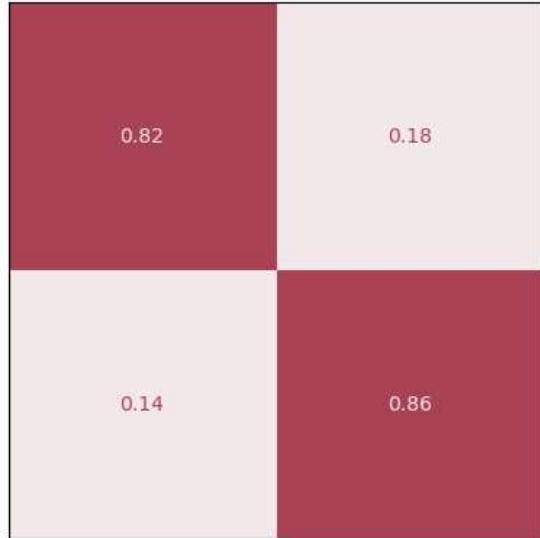
```
In [350]: # Recording and displaying the confusion matrix
best_confusion_matrix = confusion_matrix_display(best_model, y_test, best_y_pred)
```



```
In [351]: # Recording the best confusion matrix for the presentation, with labels in white
prez_cfn_matrix = best_confusion_matrix_only(best_model, y_test, best_y_pred)
```

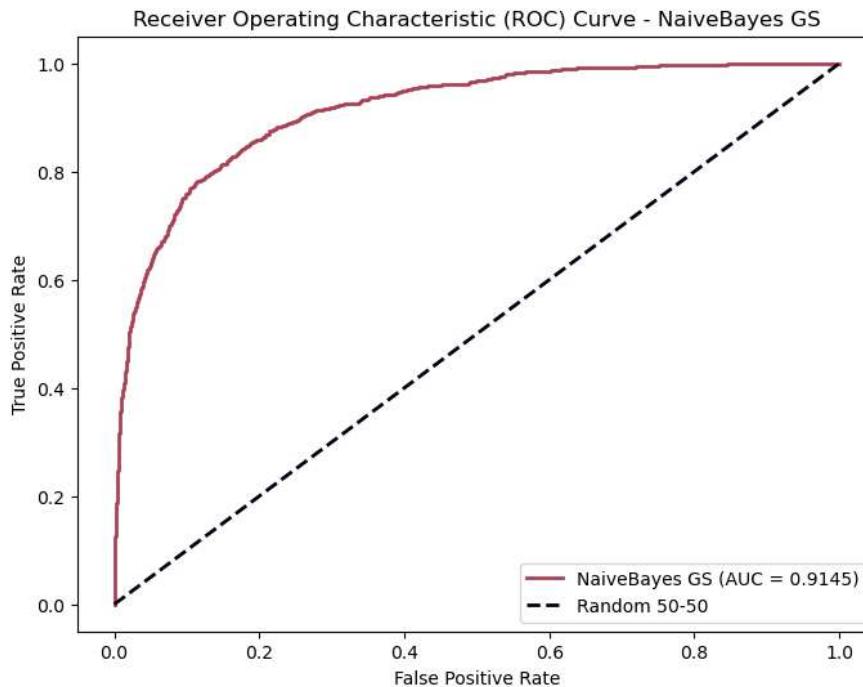


```
In [351]: # Recording the best confusion matrix for the presentation, with Labels in white
prez_cfn_matrix = best_confusion_matrix_only(best_model, y_test, best_y_pred)
```



## 5) ROC-AUC

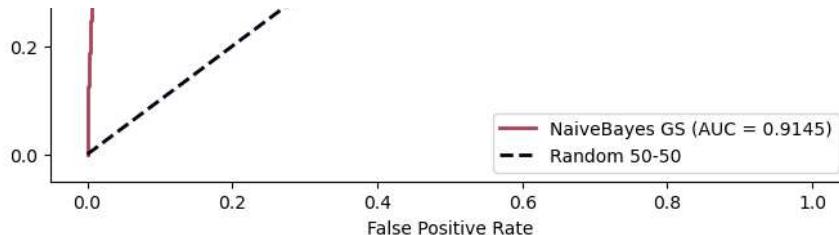
```
In [352]: # Plotting the best model ROC AUC
auc_gs = plot_roc_curve(y_test, gs_y_pred_proba, model_name=gs_model_name)
```



The balance between precision and recall. It is lower (72.5%) for the detractors class, which indicates precision is lower for this class.

The classification report, confusion matrix and ROC-AUC summarize the evaluation of the model's performance on predicting detractors from hotel TripAdvisor reviews.

The Area Under the Receiver Operating Characteristic Curve (AUC) value of 0.9145 indicates a high effectiveness in distinguishing between the two classes: `not_detractors` and `detractors`. The ROC curve, which plots the true positive rate against the false positive rate, provides a graphical representation for the evaluation of the models, as the cost of false negative is higher than the cost of false positive. If a review from a detractor was incorrectly identified as "not detractor", the hotel would miss an opportunity to transform a guest's negative experience.



The balance between precision and recall. It is lower (72.5%) for the detractors class, which indicates precision is lower for this class.  
The classification report, confusion matrix and ROC-AUC summarize the evaluation of the model's performance on predicting detractors from hotel TripAdvisor reviews.

**AUC**  
The Area Under the Receiver Operating Characteristic Curve (AUC) value of 0.9145 indicates a high effectiveness in distinguishing between the two classes: `not_detractors` and `detractors`. The ROC curve, which plots the true positive rate against the false positive rate, shows the trade-off between these two metrics. This score is higher than the random classifier, indicating better performance. The cost of false negative is higher than the cost of false positive. If a review from a detractor was incorrectly identified as "not detractor", the hotel would miss an opportunity to transform a guest's negative experience.

## 7. Findings & Recommendations

### 7. a) Most Important Features

```
In [353]: # Defining a colormap specifically for the top features bar graph
custom_colors = ['#009B8D', '#2F4858']

n_bins = 200

# Creating the custom colormap
features_cmap = LinearSegmentedColormap.from_list("features_cmap", custom_colors,
```

```
In [354]: from matplotlib.cm import ScalarMappable

def plot_top_feature_importances(pipeline, top_n):
    # Accessing the classifier instance from the pipeline
    classifier = pipeline.named_steps['classifier']

    # Getting the feature names from the vectorizer
    feature_names = pipeline.named_steps['tfidf'].get_feature_names_out()

    # Getting the feature importances
    feature_importances = classifier.feature_importances_

    # Sorting feature importances in descending order and get the corresponding indices
    sorted_indices = np.argsort(feature_importances)[-1:-top_n:-1]

    # Selecting the top 'top_n' features
    top_indices = sorted_indices[:top_n]

    # Sorting the feature names based on the selected top indices
    top_feature_names = feature_names[top_indices]
```

```
In [354]: from matplotlib.cm import ScalarMappable

def plot_top_feature_importances(pipeline, top_n):
    # Accessing the classifier instance from the pipeline
    classifier = pipeline.named_steps['classifier']

    # Getting the feature names from the vectorizer
    feature_names = pipeline.named_steps['tfidf'].get_feature_names_out()

    # Getting the feature importances
    feature_importances = classifier.feature_importances_

    # Sorting feature importances in descending order and get the corresponding indices
    sorted_indices = np.argsort(feature_importances)[::-1]

    # Selecting the top 'top_n' features
    top_indices = sorted_indices[:top_n]

    # Sorting the feature names based on the selected top indices
    top_feature_names = feature_names[top_indices]

    # Sorting the feature importances for the selected top features
    top_feature_importances = feature_importances[top_indices]

    # Reverse the order to make it descending
    top_feature_names = top_feature_names[::-1]
    top_feature_importances = top_feature_importances[::-1]

    # Creating a colormap
    cmap = plt.get_cmap(features_cmap)

    # Manually normalizing importance values
    normalized_importances = (top_feature_importances - np.min(top_feature_importances)) / (np.max(top_feature_importances) - np.min(top_feature_importances))

    # Creating the figure and axis
    fig, ax = plt.subplots(figsize=(12, 8))

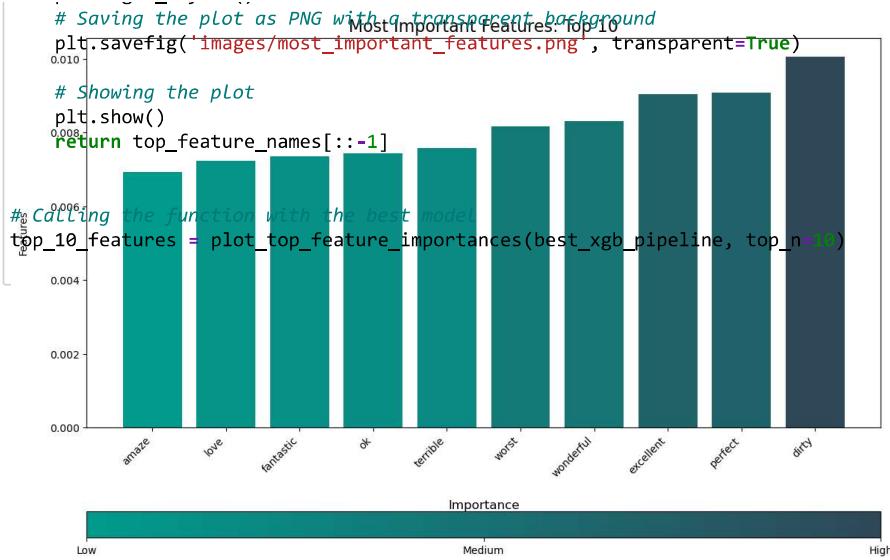
    # Using normalized_importances directly for color mapping
    bars = ax.bar(range(top_n), top_feature_importances, align='center', color=cmap(normalized_importances))
    plt.xticks(range(top_n), top_feature_names, rotation=45, ha='right') # Rotate x-axis labels
    plt.ylabel('Features', labelpad=20)
    plt.xlabel('Importance', fontsize=12, labelpad=20)
    plt.title('Most Important Features: Top {}'.format(top_n), fontsize=16)

    # Adding a colorbar below the plot
    sm = ScalarMappable(cmap=cmap)
    sm.set_array([])
    colorbar = plt.colorbar(sm, ax=ax, orientation='horizontal', aspect=30)
    # Hide colorbar ticks
    # Defining custom tick positions
    custom_ticks = [0, 0.5, 1]
    custom_ticklabels = ['Low', 'Medium', 'High']
    colorbar.set_ticks(custom_ticks)
    colorbar.set_ticklabels(custom_ticklabels)
    plt.tight_layout()
    # Saving the plot as PNG with a transparent background
    plt.savefig('images/most_important_features.png', transparent=True)

    # Showing the plot
    plt.show()
    return top_feature_names[::-1]

# Calling the function with the best model
top_10_features = plot_top_feature_importances(best_xgb_pipeline, top_n=10)
```

Feature	Importance
amaze	~0.005
love	~0.005
fantastic	~0.005
ok	~0.005
terrible	~0.006
worst	~0.007
wonderful	~0.008
excellent	~0.009
perfect	~0.009
dirty	~0.010



The most important features for the model from the 2nd best model can help hotels knowing what words to pay attention to, also when they communicate with guests. These features refer to both classes.

In [355]: # Inspecting the returned top 10 features  
top\_10\_features

Out[355]: array(['dirty', 'perfect', 'excellent', 'wonderful', 'worst', 'terrible', 'ok', 'fantastic', 'love', 'amaze'], dtype=object)

## 7. b) Recommendations

The advantage of our new service is the deployment can be immediate. Based on historical guest complaints, below are the areas we recommend to focus on to immediately start improving guest satisfactions.

1. Focus on resort hotels
  - These hotels have the most unsatisfied guests
2. Develop a maintenance program with engineering teams
  - Appearance & dysfunctionality cause frustrations
3. Train staff to enhance friendliness
  - Also ensure more languages are spoken
4. Respond to reviews
  - Detractors advise to read reviews

These areas of focus are based on the most common bigrams from detractors.

In [356]: # Plotting the detractor bigrams as a bar chart

```

# Extracting bigrams and values
bigrams, values = zip(*detractors_bigrams_10)

# Sort bigrams and values in descending order based on values
sorted_indices = sorted(range(len(values)), key=lambda k: values[k], reverse=True)
sorted_bigrams = [bigrams[i] for i in sorted_indices]
sorted_values = [values[i] for i in sorted_indices]

# Create a bar chart
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')
ax.bar(range(len(sorted_values)), sorted_values, align='center', color='#AA4255')
ax.set_xticks(range(len(sorted_values)))
ax.set_xticklabels([f'{bigram[0]} {bigram[1]}' for bigram in sorted_bigrams], rotation=45)
ax.set_xlabel('Bigrams', color='white', fontsize=13)
ax.set_ylabel('Values', color='white', fontsize=13)
ax.tick_params(axis='y', colors='white')
ax.set_title('Most Common Bigrams from Detractors', color='white', fontsize=15)

```

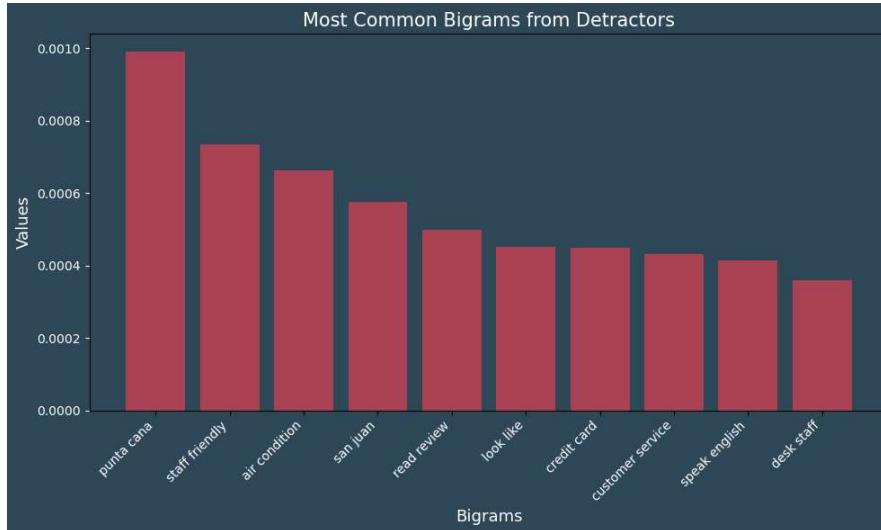
```
In [356]: # Plotting the detractor bigrams as a bar chart

# Extracting bigrams and values
bigrams, values = zip(*detractors_bigrams_10)

# Sort bigrams and values in descending order based on values
sorted_indices = sorted(range(len(values)), key=lambda k: values[k], reverse=True)
sorted_bigrams = [bigrams[i] for i in sorted_indices]
sorted_values = [values[i] for i in sorted_indices]

# Create a bar chart
fig, ax = plt.subplots(figsize=(10, 6))
fig.set_facecolor('#2F4858')
ax.set_facecolor('#2F4858')
ax.bar(range(len(sorted_values)), sorted_values, align='center', color='AA4255')
ax.set_xticks(range(len(sorted_values)))
ax.set_xticklabels([f'{bigram[0]} {bigram[1]}' for bigram in sorted_bigrams], rotation=45)
ax.set_xlabel('Bigrams', color='white', fontsize=13)
ax.set_ylabel('Values', color='white', fontsize=13)
ax.tick_params(axis='y', colors='white')
ax.set_title('Most Common Bigrams from Detractors', color='white', fontsize=15)
plt.tight_layout()

plt.savefig('images/top_10_bigrams', transparent=True)
plt.show()
```



## 8. Next Steps & Limits

- **Deployment**

Once the immediate areas of focus are implemented, the next step is only to

- entice guests to leave real-time reviews. Strategically display the QR code to facilitate

- **Next Steps**

access to the tool.  
Deploy the model for a ternary classification following the actual sentiments defined by NPS.

▪ only provide us with 2 contacts details. These persons will receive the alerts when guests score: promoters, neutral, detractors .

provide negative feedback

▪ communicate about this new tool to your teams

Very quickly, you will see the share of negative reviews decrease in your hotels, and your NPS score increase!

Extract the reviews specifically for your hotel chains, to ensure the model is based on your specific reviews. We, at TripAdvisor can do this for you for a supplement.

- **Limits**

The model provided very high scores for recommendations. If basing the model on your own hotel chains' reviews is not an option, ensure the most recent data for reviews is taken into account.

- **Next Steps**
  - Encourage guests to leave real-time reviews. Strategically display the QR code to facilitate access to the tool.
  - Deploy the model for a ternary classification, following the actual sentiments defined by NPS:
    - only provide us with 2 contacts details. These persons will receive the alerts when guests score: promoters, neutral, detractors .
    - provide negative feedback
  - communicate about this new tool to your teams

Very quickly, you will see the share of negative reviews decrease in your hotels, and your NPS score increase!

Extract the reviews specifically for your hotel chains, to ensure the model is based on your specific reviews. We, at TripAdvisor can do this for you for a supplement.

- **Limits**

The model provided very high scores for recommendations. If basing the model on your own hotel chains' reviews is not an option, ensure the most recent data for reviews is taken into account.