

```
In [ ]: ###necessary libraries
from textblob import TextBlob
import pandas as pd
import glob
import os
from datetime import datetime
import re

# file where csv files lies
path = r'C:\Users\victo\Master_Thesis\scrapersproject\bmw\bmw_scraper\spiders\news'
all_files = glob.glob(os.path.join(path, "*.csv"))

# read files to pandas frame
list_of_files = []

for filename in all_files:
    list_of_files.append(pd.read_csv(filename,
                                     sep=',',
                                     encoding='cp1252',
                                     header=None,
                                     names=["url", "header", "release time", "article content"]
                                     )
    )

# Concatenate all content of files into one DataFrames
concatenate_list_of_files = pd.concat(list_of_files,
                                       ignore_index=True,
                                       axis=0,
                                       )

# removing duplicates
cleaned_dataframe = concatenate_list_of_files.sort_values(by='url', ascending=False)
cleaned_dataframe = cleaned_dataframe.drop_duplicates(subset=["url"], keep='first', ignore_index=True)

print(cleaned_dataframe)

##formatting date column
dates = []
times = []
regex = r'(.*)((([0-2]|0?[1-9])\(/(3[01]|([12][0-9]|0?[1-9])\(/(?:[0-9]{2})?[0-9]{2})|((Jan(uary)?|Feb(ruary)?|Mar(ch)?|Apr(il)?|May|Jun(e)?|Jul(y)?|Aug(ust)?|Sep(tember)?|Oct(ober)?|Nov(ember)?|Dec(ember)?\s+\d{1,2},\s+\d{4})))'
regex2 = r'((([0-2]|0?[1-9]):([0-5][0-9])?([AaPp][Mm]))'

for date in cleaned_dataframe['release time']:
    matches = re.finditer(regex, date)
    for m in matches:
        date = m.group()
        date_formatted = date.replace(date[:2], '')
        convert_date = datetime.strptime(date_formatted, '%B %d, %Y')
        final_date = datetime.strftime(convert_date, "%Y-%m-%d")
        print(final_date)
        dates.append(final_date)

for time in cleaned_dataframe['release time']:
    matches = re.finditer(regex2, time)
    for t in matches:
        time = t.group()
        convert_time = datetime.strptime(time, '%I:%M %p')
```

```
time_formatted = datetime.strftime(convert_time, '%H:%M:%S')
print(time_formatted)
times.append(time_formatted)

## adding modified date to data frame
cleaned_dataframe['date'] = dates
cleaned_dataframe['time'] = times
cleaned_dataframe['formatted date'] = cleaned_dataframe['date'] + str(' ') + cleaned_dataframe['time']

## dropping unnecessary columns
del cleaned_dataframe['date']
del cleaned_dataframe['time']

# Join the DataFrames
cleaned_dataframe[['polarity_textblob_sentiment_header', 'subjectivity_textblob_sentiment_header']] = cleaned_dataframe['header'].apply(lambda header: pd.Series(TextBlob(header).sentiment))
cleaned_dataframe[['polarity_textblob_sentiment_content', 'subjectivity_textblob_sentiment_content']] = cleaned_dataframe['article content'].apply(lambda content: pd.Series(TextBlob(content).sentiment))

## saving outcome of flair to csv
current_date = datetime.today().strftime('%Y-%m-%d')
cleaned_dataframe.to_csv(r'C:\Users\victo\Master_Thesis\semanticanalysis\analysis_with_textblob\bmw\outcome_using_textblob\outcome_of_textblob_on_bmw_news_' + str(current_date) + '.csv', index=False)
```