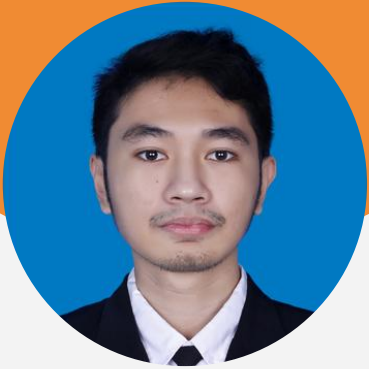# Sentiment Analysis

by : Ivan Manuel Wicaksono

# Ivan Manuel W.

**Education**

*S1 Teknik Elektro - Institut Teknologi Sepuluh Nopember (2018–2022)*

**Working**

*Business Analyst*

## Overview Project

- **Sentiment Analysis**
  Mencari trend sentimen yang terjadi melalui kanal berita website

Project Background

This project is to develop an automated web scraping from financial news (CNN and New York Times) to do a sentiment analysis (positive or negative) for each article to indicate trends

The project can help traders or investor to analyze US news sentiment (is the news give negative sentiment or positive) to decide or manage their investment portfolio

# Problem Statement

In the financial market, there is a vast amount of data available, including stock exchange information, economic indicators, global events, and financial news.

Investors can use sentiment analysis to make informed investment decisions.

The rise and fall of financial asset prices are heavily influenced by market sentiment, political news, policies, inflation data, interest rates, bond yields, and more.

To gather this information, investors need to search for news across multiple websites and look up stock and commodity prices, as this data is often scattered across the internet.

Traditionally, analysis relies on historical data, which causes delays in report generation as it takes a long time to manually read the news one by one.

With a data engineering pipeline, large market data can be processed, sentiment toward a stock can be identified, and valuable insights can be generated.

The dashboard displays:

- Sentiment analysis for each article
- Comparison with SP500

# Data Platform Understanding

# Timeline

**Data Extraction**

Scrape financial news website

**Data Processing**

Cleaning null values and duplicates values with spark

**Data Storage**

Load the data to Bigquery

**Data Visualization**

Visualize data in Tableau

ibimbing

# Data Understanding

Datasets :

- Websites (CNN, New York Times 250 articles in the last 7 days)
- CSV (SP500 price in the last 7 days)

Data Orchestration :

- Using airflow for scheduling and monitoring batch processing

Extract :

- Web scraping the data with python and selenium

Transform :

- Clean the data, duplicate values, and missing values using pyspark
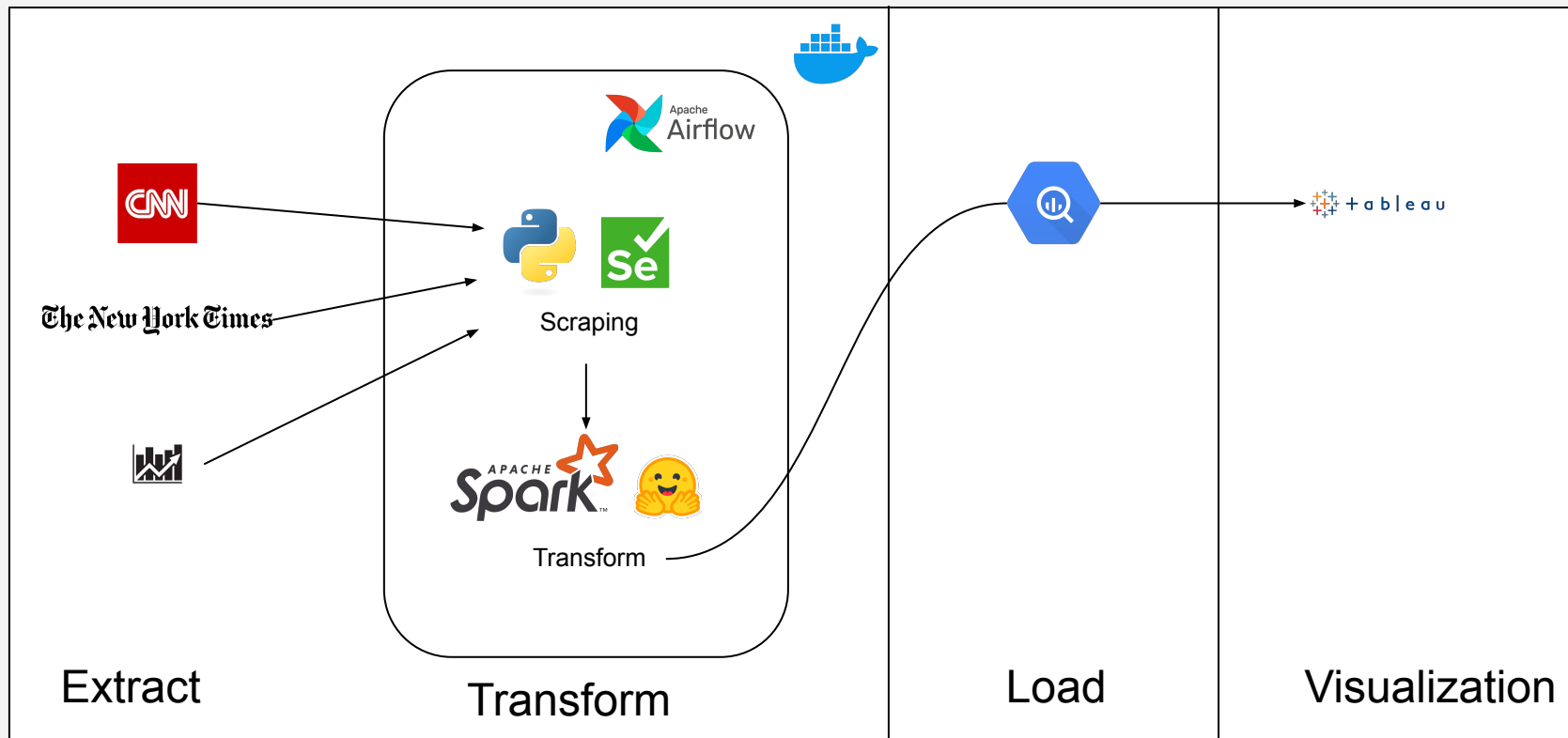- Do a sentiment analysis with ML model hugging face

Load :

- Load the data to BigQuery

# Transformation & Consideration

# Architecture



Extract · Transform · Load · Visualization

| CNN | |
|---|---|
| **PK** | **datetime** |
| | close |
| | open |
| | high |
| | low |

| CNN | |
|---|---|
| **PK** | **url** |
| | datetime |
| | title |
| | text |
| | source |
| | label |
| | score |

| NYT | |
|---|---|
| **PK** | **url** |
| | datetime |
| | title |
| | text |
| | source |
| | label |
| | score |

Conclusion & Recommendation

# DAG Airflow - News

# DAG Airflow - SP500

# Load BigQuery

# Perform Query

# Code Overview - Web Scraping

ibimbing

```
kefile        scraping-news-dag.py      nyt.py  1 ×      cnn.py      ...        cnn.py ×                                         ▷ ∨   ⬚   ...

dags > source >   nyt.py >  transform                    dags > source >   cnn.py >  get_all_url
   1   from selenium import webdriver                        1   def get_all_url(base_url):
   2   from selenium.webdriver.chrome.options import Optio    2       import requests
   3   from selenium.webdriver.common.by import By            3       from bs4 import BeautifulSoup
   4   from selenium.webdriver.support.ui import WebDriver    4
   5   from selenium.webdriver.support import expected_con    5       res = requests.get(base_url)
   6   from selenium.webdriver.chrome.service import Servi    6       soup = BeautifulSoup(res.text, 'html.parser')
   7   from webdriver_manager.chrome import ChromeDriverMa    7
   8   from bs4 import BeautifulSoup                           8       list_content = soup.find_all('a', {'class': 'co
   9   import newspaper                                        9       http_url = 'https://edition.cnn.com'
  10   import pandas as pd                                    10
  11   from datetime import datetime, timedelta              11       url_date = []
  12                                                          12       for content in list_content:
  13   def setup_driver():                                   13           data = dict()
  14       # Initialize webdriver                            14           url = f"{http_url}{content['href']}"
  15       options = Options()                               15           url_date.append(url)
  16       options.add_argument("--headless")  # Run Chrom   16
  17       options.add_argument("--no-sandbox")             17       filtered_url = set(url_date)
  18       options.add_argument("--disable-dev-shm-usage")  18       return filtered_url
  19       driver = webdriver.Chrome(service=Service(Chrom  19
  20       return driver                                     20   def transform(url_article):
  21                                                          21       # Transform datetime EST to UTC + 7
  22   def get_page_source(driver, url):                     22       import requests
  23       # get page_source                                 23       from bs4 import BeautifulSoup
  24       driver.get(url)                                   24       from datetime import datetime
  25       return driver.page_source                         25
  26                                                          26       res = requests.get(url_article)
```
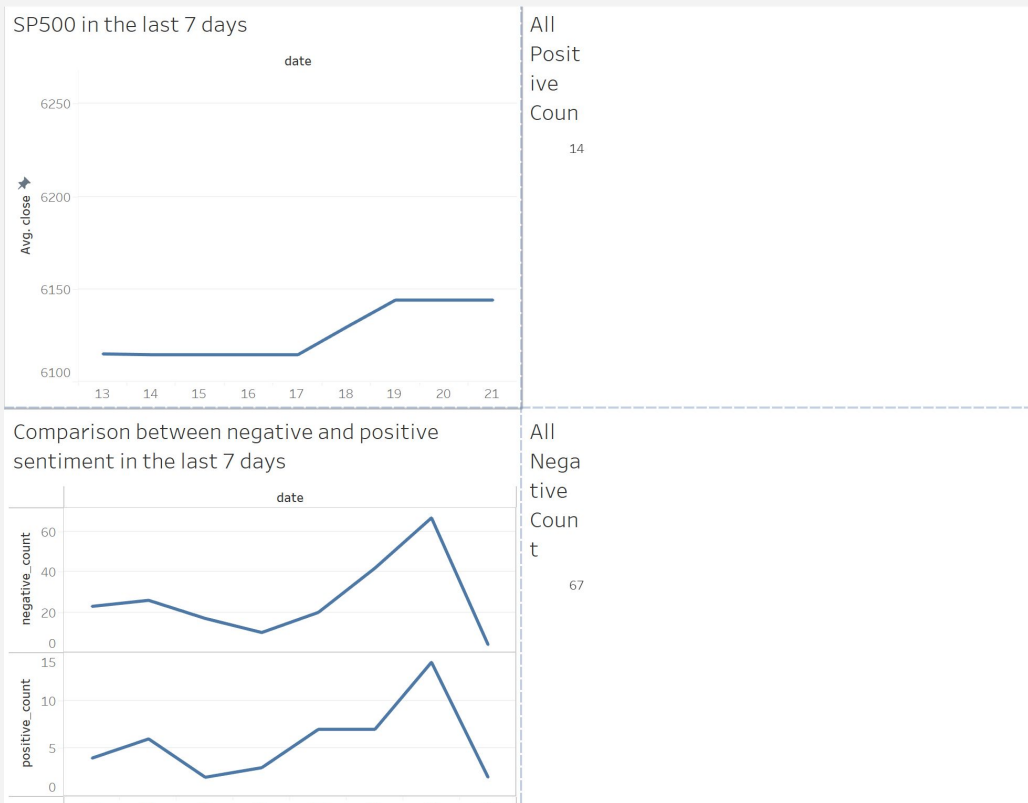
These are the codes for web scraping through CNN and NYT. Because NYT is a dynamic website, I use selenium to scrap dynamic website

# Code Overview - DAG



```python
1   from airflow.decorators import dag, task
2   from airflow.operators.empty import EmptyOperator
3   from airflow.providers.apache.spark.operators.spark_submit import SparkSubmitOperator
4   from google.oauth2 import service_account
5   import yaml
6
7   with open("dags/source/list_tables.yaml") as f:
8       list_tables = yaml.safe_load(f)
9
10  @dag()
11  def scraping_news_dag():
12      start_task = EmptyOperator(task_id="start_task")
13      end_task = EmptyOperator(task_id="end_task")
14
15      for table in list_tables:
16          spark_submit = SparkSubmitOperator(
17              application = f"/spark-scripts/spark-{table}.py",
18              conn_id = "spark_main",
19              task_id = f"spark_load_task_{table}",
20          )
21
22          @task(task_id=f"scraping_{table}")
23          def scrapping(table_name):
24              from source import nyt, cnn
25              module_map = {"cnn": cnn, "nyt": nyt}
26              df = module_map[table_name].main()
```

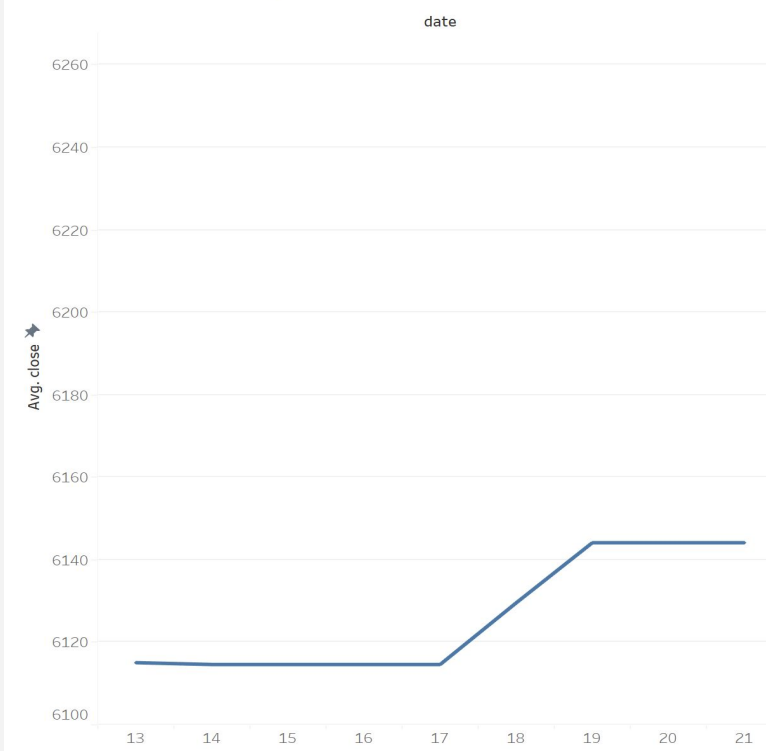This is dynamic dag to scrap multiple website. You can add more if you want to

ibimbing

# Visualization



SP500 in the last 7 days

All Positive Count

Comparison between negative and positive sentiment in the last 7 days

All Negative Count

# Visualization

# Conclusion

We can see from the dashboard that negative sentiment dominate the news. The ratio between positive and negative news are 0.1 to 0.3. So, we can conclude that the news in last 7 days are heavily negative. But if we compare with sp500, the index has increased from 6115 to 6140 or 0.4% gains at 17 February until 19 February.

| Date | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|------|------|------|------|------|------|------|------|------|
| Positive | 4 | 6 | 2 | 3 | 7 | 7 | 14 | 2 |
| Negative | 23 | 26 | 17 | 10 | 20 | 42 | 67 | 4 |
| Ratio | 0.174 | 0.231 | 0.118 | 0.3 | 0.35 | 0.167 | 0.209 | 0.5 |

# Conclusion

From this analysis we can conclude that we can't predict the stock price using only one metric. This project only one of supporting tools to help analyze US sentiment

The limitation of this project that it's not sufficient using only 2 website sources. More data can help to generalize the news and get broader insight. But it is hard to find news website that can be scrape. Because I am using built-in model machine learning library to give sentiment analysis, there is a limitation for tokenize words.

**Terima Kasih.**

ibimbing