資工碩一 113598007 楊明哲

環境設置

1. Windows

• 環境版本

Java: 17.0.12

Python: 3.10.10

• PySpark: 3.5.3

Pandas: 2.2.3

- 執行指令 spark-class org.apache.spark.deploy.master.Master,並打開 URL
 - 指令結果

```
65 D:Code\Big_Data Mining> spark-class org.apache.spark.deploy.master.Master
Using Spark's default logdj profile: org/apache/spark/logdj2-defaults.properties
24/10/16 23:33:15 IMFO Master: Started daemon with process name: 121480DESKTOP-F8USCHM
24/10/16 23:33:15 IMFO Master: Started daemon with process name: 121480DESKTOP-F8USCHM
24/10/16 23:33:15 IMFO SecurityManager: Changing wiew acls to: Zhe
24/10/16 23:33:16 IMFO SecurityManager: Changing modify acls to: Zhe
24/10/16 23:33:16 IMFO SecurityManager: Changing modify acls to: Zhe
24/10/16 23:33:16 IMFO SecurityManager: Changing woidify acls groups to:
24/10/16 23:33:16 IMFO SecurityManager: Changing modify acls groups to:
24/10/16 23:33:16 IMFO SecurityManager: SecurityManager according to disabled; ui acls disabled; users with view permissions: Zhe; groups with view permissions: EMPTY; users with vodify permissions: Zhe; groups with modify permissions: EMPTY; users with vodify permissions: Zhe; groups with modify permissions: EMPTY
24/10/16 23:33:16 IMFO Master: Starting Spark master at spark://192.168.18.5:7077
24/10/16 23:33:16 IMFO Master: Starting Spark master at spark://192.168.18.5:7077
24/10/16 23:33:16 IMFO Master: Starting Spark master at spark://192.168.18.5:7077
24/10/16 23:33:16 IMFO Master: Starting Spark waster at spark://192.168.18.5:7077
24/10/16 23:33:16 IMFO Master: Linuming Spark version 3.5.3
24/10/16 23:33:16 IMFO Master: Linuming Spark version 3.5.3
24/10/16 23:33:16 IMFO Master: Implied Spark waster at spark://192.168.18.5:7077
24/10/16 23:33:16 IMFO Master: Registering worker 192.168.46.128:44857 with 16 cores, 6.7 GiB RAM
```

2. Ubuntu

• 環境版本

Memory: 8GB

• Processors: 16

Java: 17.0.12

• Python: 3.10.12

• Spark: 3.5.3

• 執行 spark-class org.apache.spark.deply.worker.Worker spark://192.168.18.5:7077

資料夾結構

- homework1.py
- archive
 - spacenews.csv
- image
 - all of output snapshots

源代碼

```
import pandas as pd
 from pyspark.sql import SparkSession
 from pyspark.sql.functions import explode, split, col, lower, regexp_replace, count, to_date, date_format
# Create SparkSession
spark = SparkSession.builder.appName("Word Frequency").getOrCreate()
# Read Data
file_path = './archive/spacenews.csv'
df = pd.read_csv(file_path)
df.dropna(inplace = True)
spark_df_origin = spark.createDataFrame(df)
 spark_df_origin = spark_df_origin.withColumn('Date', to_date(col('Date'), 'MMMM d, yyyy'))
spark_df = spark_df_origin
spark_df = spark_df.withColumn('Title', lower(col('Title')))
spark_df = spark_df.withColumn('Title', regexp_replace(col('Title'), '[^\w\s]', ''))
spark_df = spark_df.withColumn('word', explode(split(col('Title'), '\s+')))
 total_word_freq = spark_df.groupBy('word').count().orderBy(col('count').desc())
 total_word_freq = spark_df.groupBy('word').count().orderBy(col('count').desc())
 daily_word_freq = spark_df.groupBy('Date', 'word').count().orderBy(col('Date'), col('count').desc())
 total_word_freq.show()
 daily_word_freq.show()
 spark_df = spark_df_origin
 spark_df = spark_df.withColumn('Content', lower(col('Content')))
 spark\_df = spark\_df.withColumn('Content', regexp\_replace(col('Content'), '[^\setminus w \setminus s]', ''))
 spark_df = spark_df.withColumn('content_word', explode(split(col('Content'), '\s+')))
 total_content_word_freq = spark_df.groupBy('content_word').count().orderBy(col('count').desc())
 daily_content_word_freq = spark_df.groupBy('Date', 'content_word').count().orderBy(col('Date'), col('count').desc())
 total_content_word_freq.show()
daily_content_word_freq.show()
 # Q3 : Daily article count
 spark_df = spark_df_origin
daily_article_count = spark_df.groupBy('Date').count().withColumnRenamed('count', 'daily_count')
daily_article_percentage = daily_article_count.withColumn('percentage', (col('daily_count') / total_article_count) * 100)
author_daily_article_count = spark_df.groupBy('Date', 'Author').count().withColumnRenamed('count', 'author_daily_count')
author_daily_article_percentage = author_daily_article_count.join(daily_article_count, 'Date')
author_daily_article_percentage = author_daily_article_percentage.withColumn('percentage', (col('author_daily_count')) / col('daily_count')) * 100)
daily_article_percentage.show()
author_daily_article_percentage.show()
spark df = spark df origin
spark_df = spark_df.withColumn('Title', lower(col('Title')))
spark_df = spark_df.withColumn('Postexcerpt', lower(col('Postexcerpt')))
spark_df = spark_df.withColumn('Title', regexp_replace(col('Title'), '[^\w\s]', ''))
spark_df = spark_df.withColumn[]'Postexcerpt', regexp_replace(col('Postexcerpt'), '[^\w\s]', '')[]
filtered_df = spark_df.filter((col('Title').contains('space'))) & (col('Postexcerpt').contains('space')))
sorted_filtered_df = filtered_df.orderBy('Date')
sorted filtered df.show()
```

輸出結果

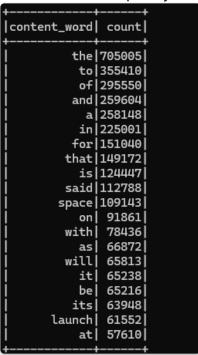
執行 spark-submit --master spark://192.168.18.5:7077 --conf spark.driver.host=192.168.18.5 homework1.py

• Title word frequency

	1 7
word	count
to	5786
space	4113
for	3980
of	2254
in	2208
launch	2084
satellite	2024
on	1922
nasa	1887
the	1427
force	1286
l us	1249
and and	1242
with	1090
a	997
new	984
satellites	901
commercial	851
spacex	837
air	780
+	+

·		
+ Date		+ count
Date	WOIG	Counci
2005-08-01	in	1
12005-08-01	:	: :
2005-08-01		
2005-08-01		
2005-08-01		
2005-08-01		
2005-08-01		
2005-08-01		
2005-11-07		
2005-11-07		
2005-11-07		: :
2005-11-07		-!
2005-11-07		
2005-11-07		: :
2005-11-07		
2005-11-07		: :
2005-11-07		: :
:	commercial	
2005-11-21		: :
2005-11-21	: -	: :
+		

Content word frequency



Date	 content_word	count
2005-08-01	the	45
2005-08-01	of	18
2005-08-01	galileo	17
2005-08-01	to	16
2005-08-01	is	16
2005-08-01	in	15
2005-08-01	that	14
2005-08-01	said	12
2005-08-01	and	12
2005-08-01	on on	10
2005-08-01	will	10
2005-08-01	a	10
2005-08-01	be	10
2005-08-01	with	9
2005-08-01	esa	8
2005-08-01	by	8
2005-08-01	we	7
2005-08-01	funding	7
2005-08-01		
2005-08-01	million	7
+	+	

第三題

Daily article count

<u> </u>		
Date	daily_count	percentage
2022-07-31	3	0.01663708961845608
2022-11-29	10	0.05545696539485359
2022-07-27	8	0.044365572315882874
2022-08-02	5	0.027728482697426796
2022-05-19	7	0.03881987577639751
2022-05-26	6	0.03327417923691216
2022-06-22	9	0.04991126885536823
2022-06-06	3	0.01663708961845608
2022-11-21	4	0.022182786157941437
2022-07-07	7	0.03881987577639751
2022-12-09	6	0.03327417923691216
2022-12-02	4	0.022182786157941437
2022-10-19	5	0.027728482697426796
2022-08-11	8	0.044365572315882874
2022-08-20	1	0.005545696539485359
2022-09-29	8	0.044365572315882874
2022-11-20	2	0.011091393078970719
2022-12-11	3	0.01663708961845608
2022-09-06	6	0.03327417923691216
2022-09-07	8	0.044365572315882874
(·+

Date	Author	author_daily_count	daily_count	percentage
2022-07-31	Jeff Foust			180.0
2022-11-29	Sandra Erwin		10	20.0
2022-11-29	Jeff Foust		10	30.0
2022-11-29	Debra Werner			10.0
2022-11-29	Park Si-soo			10.0
2022-11-29	Andrew Jones		10	10.0
2022-11-29	Jason Rainbow			28.8
2022-07-27	Andrew Jones			12.5
2022-07-27	Debra Werner			12.5
2022-07-27	Sandra Erwin			25.0
2022-07-27	Jeff Foust			25.8
2022-07-27	Jason Rainbow			25.0
	Jeff Foust			40.0
2022-08-02	Park Si-soo			28.8
2022-08-02	Jason Rainbow			20.0
2022-08-02	Sandra Erwin			20.0
	Sandra Erwin			14.285714285714285
	Jason Rainbow			28.57142857142857
2022-05-19	Jeff Foust			57.14285714285714
2022-05-26	Sandra Erwin			33.3333333333333

第四題

• Space in Title and Postexcerpt

	url	content	author	Date	Postexcerpt
despite the budge lightfoot to lead mike gold corpora peckham leaves sp spacex lands orbc ariane 5 remnant spacex hardware d amid shift in pow japans first spac spacex and astriu white house hande us canada sign sp arianespace chief esa designing doc us france sign fo a 50year history volker liebig dir japans htvl cargo spacex says drago minotaur 4 concer	https://spacenews	WASHINGTON — As t WASHINGTON — NASA WITH The administ WASHINGTON — Rob WASHINGTON — Rob WASHINGTON — Sate The remains of a Space Exploration TOKYO — The Japan NEW YORK — Japan WASHINGTON — Hawt WASHINGTON — Add NASA Administrato PARIS — Engineers PARIS — The Europ NASA Administrato BERKELEY, Calif Satellite Earth o Japan's first H-2 WASHINGTON — Spac	Brian Berger Amy Klamper Amy Klamper Amy Klamper Warren Ferster SpaceNews Staff SpaceNews Staff Paul Kallender-Umezu SpaceNews Staff Amy Klamper Amy Klamper SpaceNews Staff Peter B. de Selding Peter B. de Selding SpaceNews Staff Debra Werner Peter B. de Selding SpaceNews Staff Debra Werner Peter B. de Selding	2006-08-29 2009-08-21 2009-08-31 2009-09-03 2009-09-07 2009-09-07 2009-09-07 2009-09-07 2009-09-08 2009-09-14 2009-09-14 2009-09-18 2009-09-18 2009-09-21 2009-09-21 2009-09-21	as the new direct washington nasa with the administ washington rob p washington satel the remains of a space exploration

小組分工

學號、姓名	貢獻比例	工作內容
楊明哲	100%	程式編寫、除錯以及文書處理