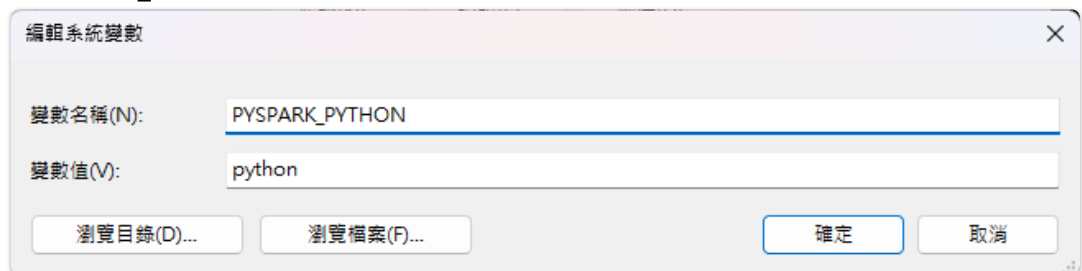


# 資工碩一 113598007 楊明哲

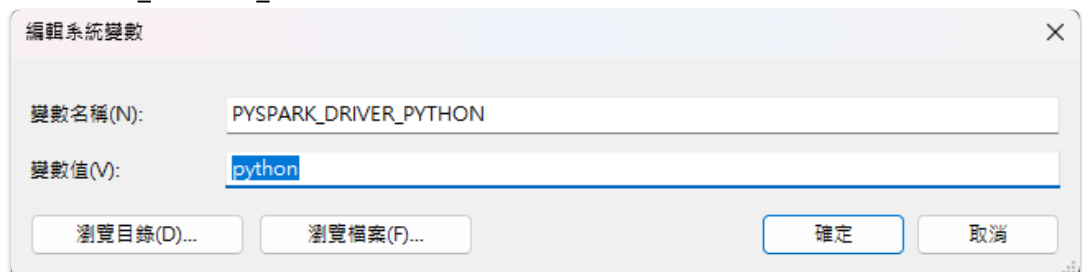
## 環境設置

### 1. Windows

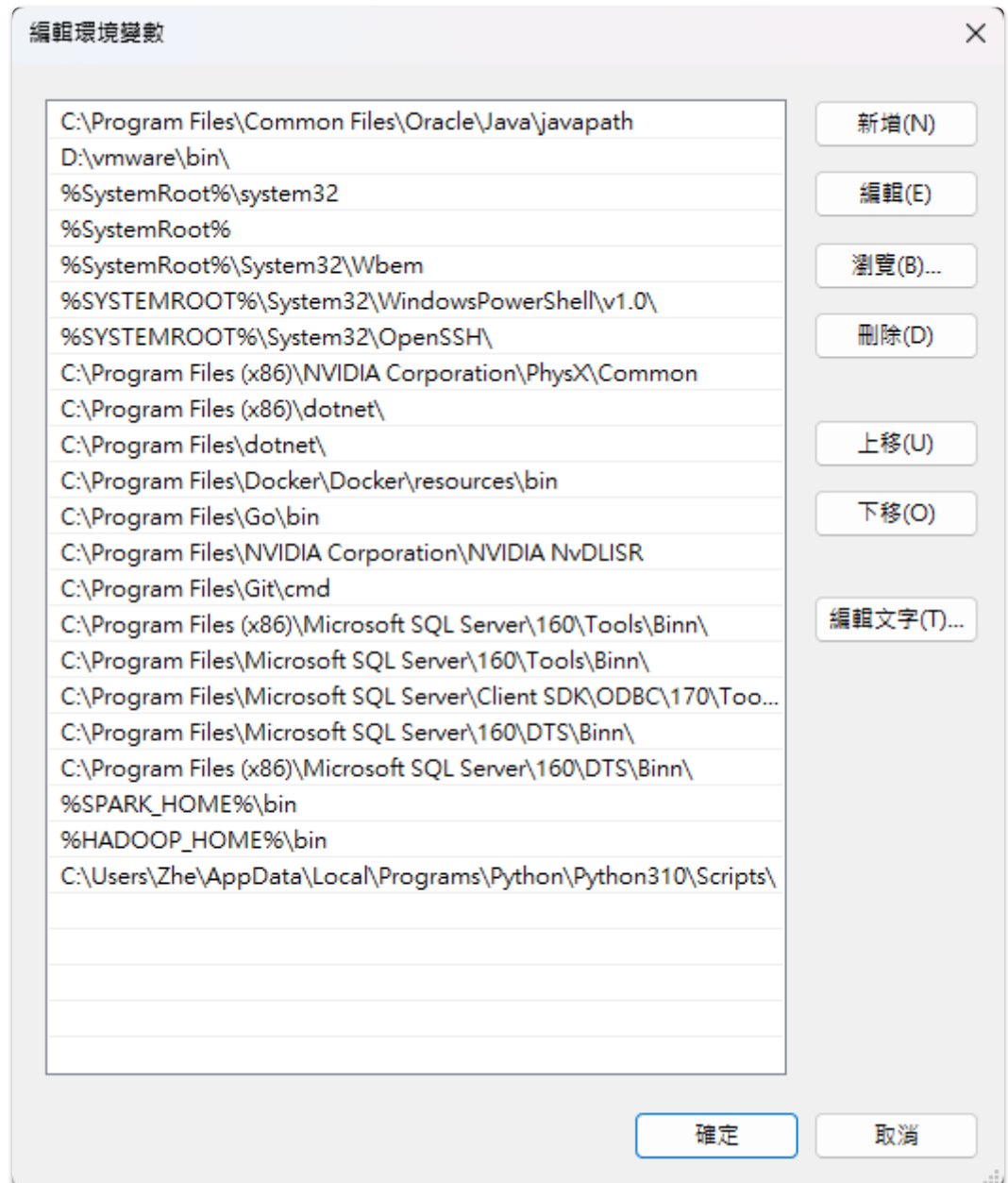
- 環境版本
  - Java : 17.0.12
  - Python : 3.10.10
  - PySpark : 3.5.3
  - Pandas : 2.2.3
- 環境配置
  - 下載指定版本的 Python 和 Java
  - 透過 `pip install pyspark` 安裝 pyspark，並且將它初始化
  - 將下載的 `winutils.exe` 放進 pyspark 中 `bin` 的資料夾下
  - 設置環境變數
    - PYSPARK\_PYTHON



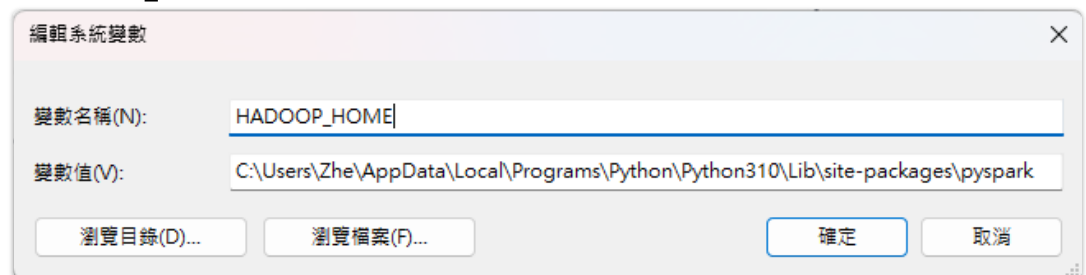
- PYSPARK\_DRIVER\_PYTHON



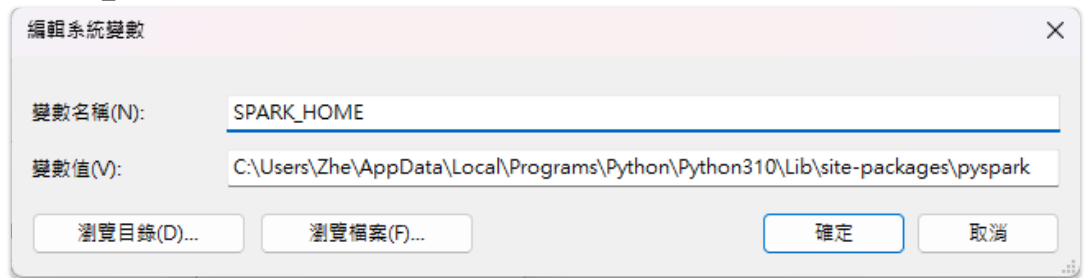
- Path



- HADOOP\_HOME



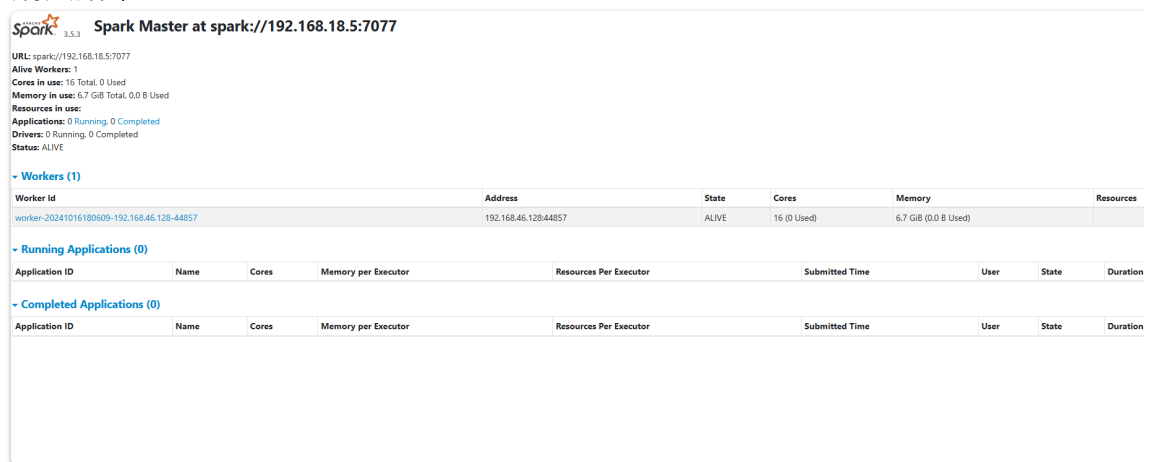
- SPARK\_HOME



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

```
PS D:\Code\Big_Data_Mining> spark-class org.apache.spark.deploy.master.Master
Using Spark's default log4j profile: org/apache/spark/log4j2-defaults.properties
24/10/16 23:33:15 INFO Master: Started daemon with process name: 12148@DESKTOP-F8USCHM
24/10/16 23:33:15 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
24/10/16 23:33:16 INFO SecurityManager: Changing view acls to: Zhe
24/10/16 23:33:16 INFO SecurityManager: Changing modify acls to: Zhe
24/10/16 23:33:16 INFO SecurityManager: Changing view acls groups to:
24/10/16 23:33:16 INFO SecurityManager: Changing modify acls groups to:
24/10/16 23:33:16 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: Zhe; groups with view permissions: EMPTY; u
odify permissions: Zhe; groups with modify permissions: EMPTY
24/10/16 23:33:16 INFO Utils: Successfully started service 'sparkMaster' on port 7077.
24/10/16 23:33:16 INFO Master: Starting Spark master at spark://192.168.18.5:7077
24/10/16 23:33:16 INFO Master: Running Spark version 3.5.3
24/10/16 23:33:16 INFO JettyUtils: Start Jetty 0.0.0.0:8080 for MasterUI
24/10/16 23:33:16 INFO Utils: Successfully started service 'MasterUI' on port 8080.
24/10/16 23:33:16 INFO MasterWebUI: Bound MasterWebUI to 0.0.0.0, and started at http://DESKTOP-F8USCHM:8080
24/10/16 23:33:16 INFO Master: I have been elected leader! New state: ALIVE
24/10/16 23:33:18 INFO Master: Registering worker 192.168.46.128:44857 with 16 cores, 6.7 GiB RAM
```

- 網頁結果



## 2. Ubuntu

- 環境版本
  - Memory : 8GB
  - Processors : 16
  - Java : 17.0.12
  - Python : 3.10.12
  - Spark : 3.5.3
- 環境配置
  - 下載指定的 Java 版本
  - 下載指定的 Spark，並且解壓縮到 Documents 下

- 將各環境變數設置

```
cluster0@cluster0-virtual-machine: ~  
export JAVA_HOME=/usr/lib/jvm/java-17-openjdk-amd64  
export SPARK_HOME=/home/cluster0/Documents/spark-3.5.3-bin-hadoop3  
export HADOOP_HOME=/home/cluster0/Documents/spark-3.5.3-bin-hadoop3  
export PATH=$PATH:$SPARK_HOME/bin:$HADOOP_HOME/sbin  
  
# ~/.bashrc: executed by bash(1) for non-login shells.  
# see /usr/share/doc/bash/examples/startup-files (in the package bash-doc)  
# for examples  
  
# If not running interactively, don't do anything  
case $- in  
    *(*) ;;  
    *) return;;  
esac  
  
# don't put duplicate lines or lines starting with space in the history.  
# See bash(1) for more options  
HISTCONTROL=ignoreboth  
  
# append to the history file, don't overwrite it  
shopt -s histappend
```

- 執行 `spark-class org.apache.spark.deploy.worker.Worker spark://192.168.18.5:7077`

```
cluster0@cluster0-virtual-machine:~/Documents$ spark-class org.apache.spark.deploy.worker.Worker spark://192.168.18.5:7077  
Using Spark's default log4j profile: org/apache/spark/log4j2-defaults.properties  
24/10/16 18:06:09 INFO Worker: Started daemon with process name: 19654@cluster0-virtual-machine  
24/10/16 18:06:09 INFO SignalUtils: Registering signal handler for TERM  
24/10/16 18:06:09 INFO SignalUtils: Registering signal handler for HUP  
24/10/16 18:06:09 INFO SignalUtils: Registering signal handler for INT  
24/10/16 18:06:09 WARN Utils: Your hostname, cluster0-virtual-machine resolves to a loopback address: 127.0.1.1; using 192.168.46.128 instead (on interface ens33)  
24/10/16 18:06:09 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address  
24/10/16 18:06:09 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
24/10/16 18:06:09 INFO SecurityManager: Changing view acls to: cluster0  
24/10/16 18:06:09 INFO SecurityManager: Changing modify acls to: cluster0  
24/10/16 18:06:09 INFO SecurityManager: Changing view acls groups to:  
24/10/16 18:06:09 INFO SecurityManager: Changing modify acls groups to:  
24/10/16 18:06:09 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: cluster0; groups with view permissions: EMPTY; users with modify permissions: cluster0; groups with modify permissions: EMPTY  
24/10/16 18:06:09 INFO Utils: Successfully started service 'sparkWorker' on port 44857.  
24/10/16 18:06:09 INFO Worker: Worker decommissioning not enabled.  
24/10/16 18:06:10 INFO Worker: Starting Spark worker 192.168.46.128:44857 with 16 cores, 6.7 GiB RAM  
24/10/16 18:06:10 INFO Worker: Running Spark version 3.5.3  
24/10/16 18:06:10 INFO Worker: Spark home: /home/cluster0/Documents/spark-3.5.3-bin-hadoop3  
24/10/16 18:06:10 INFO ResourceUtils: =====  
24/10/16 18:06:10 INFO ResourceUtils: No custom resources configured for spark.worker.  
24/10/16 18:06:10 INFO ResourceUtils: =====  
24/10/16 18:06:10 INFO JettyUtils: Start Jetty 0.0.0.0:8081 for WorkerUI  
24/10/16 18:06:10 INFO Utils: Successfully started service 'WorkerUI' on port 8081.  
24/10/16 18:06:10 INFO WorkerWebUI: Bound WorkerWebUI to 0.0.0.0, and started at http://192.168.46.128:8081
```

- 網頁畫面



- 並且能夠在 Master 上的網頁看到 Worker 的資訊

Workers (1)					
Worker ID	Address	State	Cores	Memory	Resources
worker-20241016180609-192.168.46.128-44857	192.168.46.128:44857	ALIVE	16 (0 Used)	6.7 GiB (0.0 B Used)	

## 源代碼

```
1 import pandas as pd
2 from pyspark.sql import SparkSession
3 from pyspark.sql.functions import col, min, max, mean, stddev
4
5 # Create SparkSession
6 spark = SparkSession.builder.appName("Electric Power Consumption").getOrCreate()
7
8 # Read Data
9 file_path = "D:/Code/Big_Data_Mining/Electric_Power_Consumption/household_power_consumption.txt"
10 df = pd.read_csv(file_path, sep=';', low_memory=False)
11
12 # Drop missing values
13 df.dropna(inplace=True)
14
15 # Data type conversion
16 df['Global_active_power'] = df['Global_active_power'].astype(float)
17 df['Global_reactive_power'] = df['Global_reactive_power'].astype(float)
18 df['Voltage'] = df['Voltage'].astype(float)
19 df['Global_intensity'] = df['Global_intensity'].astype(float)
20
21 # Convert pandas DataFrame to Spark DataFrame
22 spark_df = spark.createDataFrame(df)
23
24 # Q1: Find min, max, count
25 min_values = spark_df.select(
26     min(col('Global_active_power')).alias('Global_active_power'),
27     min(col('Global_reactive_power')).alias('Global_reactive_power'),
28     min(col('Voltage')).alias('Voltage'),
29     min(col('Global_intensity')).alias('Global_intensity')
30 ).collect()[0]
31
32 max_values = spark_df.select(
33     max(col('Global_active_power')).alias('Global_active_power'),
34     max(col('Global_reactive_power')).alias('Global_reactive_power'),
35     max(col('Voltage')).alias('Voltage'),
36     max(col('Global_intensity')).alias('Global_intensity')
37 ).collect()[0]
38
39 count_values = spark_df.count()
40
41 stats1 = pd.DataFrame({
42     'Min': [min_values['Global_active_power'], min_values['Global_reactive_power'], min_values['Voltage'], min_values['Global_intensity']],
43     'Max': [max_values['Global_active_power'], max_values['Global_reactive_power'], max_values['Voltage'], max_values['Global_intensity']],
44     'Count': [count_values, count_values, count_values, count_values]
45 }, index=['Global_active_power', 'Global_reactive_power', 'Voltage', 'Global_intensity'])
46
47 print("Q1:")
48 print(stats1)
49
50 # Q2: Find mean, standard deviation
51 mean_values = spark_df.select(
52     mean(col('Global_active_power')).alias('Global_active_power'),
53     mean(col('Global_reactive_power')).alias('Global_reactive_power'),
54     mean(col('Voltage')).alias('Voltage'),
55     mean(col('Global_intensity')).alias('Global_intensity')
56 ).collect()[0]
57
58 stddev_values = spark_df.select(
59     stddev(col('Global_active_power')).alias('Global_active_power'),
60     stddev(col('Global_reactive_power')).alias('Global_reactive_power'),
61     stddev(col('Voltage')).alias('Voltage'),
62     stddev(col('Global_intensity')).alias('Global_intensity')
63 ).collect()[0]
64
65 stats2 = pd.DataFrame({
66     'Mean': [mean_values['Global_active_power'], mean_values['Global_reactive_power'], mean_values['Voltage'], mean_values['Global_intensity']],
67     'StdDev': [stddev_values['Global_active_power'], stddev_values['Global_reactive_power'], stddev_values['Voltage'], stddev_values['Global_intensity']]
68 }, index=['Global_active_power', 'Global_reactive_power', 'Voltage', 'Global_intensity'])
69
70 print("Q2:")
71 print(stats2)
72
73 # Q3: Min-Max Normalization
74 normalized_df = spark_df.select(
75     ((col('Global_active_power') - min_values['Global_active_power']) / (max_values['Global_active_power'] - min_values['Global_active_power'])).alias('Normalized_Global_active_pow
76     ((col('Global_reactive_power') - min_values['Global_reactive_power']) / (max_values['Global_reactive_power'] - min_values['Global_reactive_power'])).alias('Normalized_Global_re
77     ((col('Voltage') - min_values['Voltage']) / (max_values['Voltage'] - min_values['Voltage'])).alias('Normalized_Voltage'),
78     ((col('Global_intensity') - min_values['Global_intensity']) / (max_values['Global_intensity'] - min_values['Global_intensity'])).alias('Normalized_Global_intensity')
79 )
80
81 # Write to CSV
82 stats1_output_path = "D:/Code/Big_Data_Mining/result/Q1.csv"
83 stats2_output_path = "D:/Code/Big_Data_Mining/result/Q2.csv"
84 normalized_output_path = "D:/Code/Big_Data_Mining/result/Q3.csv"
85
86 stats1.to_csv(stats1_output_path)
87 stats2.to_csv(stats2_output_path)
88 normalized_df.toPandas().to_csv(normalized_output_path, index=False)
```

## 輸出結果

- 執行 `spark-submit --master spark://192.168.18.5:7077 --conf spark.driver.host=192.168.18.5 homework0.py`

## 第一題

- Find min, max, count

Q1:

	Min	Max	Count
Global_active_power	0.076	11.122	2049280
Global_reactive_power	0.000	1.390	2049280
Voltage	223.200	254.150	2049280
Global_intensity	0.200	48.400	2049280

## 第二題

- Find mean, standard deviation

Q2:

	Mean	StdDev
Global_active_power	1.091615	1.057294
Global_reactive_power	0.123714	0.112722
Voltage	240.839858	3.239987
Global_intensity	4.627759	4.444396

## 第三題

- Normalization

```
1 Normalized_Global_active_power,Normalized_Global_reactive_power,Normalized_Voltage,Normalized_Global_intensity
2 0.3747963063552418,0.30071942446043165,0.376090468497577,0.37759336099585067
3 0.4783632084012313,0.31366906474820144,0.33699515347334413,0.47302904564315357
4 0.4796306355241717,0.35827338129496406,0.32600969305331173,0.47302904564315357
5 0.48089806264711216,0.3611510791366907,0.3405492730210021,0.47302904564315357
6 0.325004526524391,0.3798561151079137,0.4032310177705981,0.32365145228215775
7 0.311787072243346,0.37553956834532376,0.3819063004846531,0.3070539419087137
8 0.3282636248415716,0.3741007194244605,0.3841680129240713,0.32365145228215775
9 0.32808256382400874,0.3741007194244605,0.38836833602584825,0.32365145228215775
10 0.3251858754300203,0.36690647482014394,0.3486268174474964,0.32365145228215775
11 0.32464240449031323,0.36690647482014394,0.34442649434571954,0.32365145228215775
12 0.39579938439254037,0.35827338129496406,0.31211631663974215,0.40248962655601667
13 0.48307079485786714,0.3381294964028777,0.30953150242326355,0.47717842323651455
14 0.46605105920695283,0.34388489208633094,0.3163166397415191,0.46058091286307057
15 0.470034401593337,0.2863309352517986,0.3137318255250405,0.4647302904564316
16 0.36013036393264536,0.3035971223021583,0.3890145395799681,0.3609958506224067
17 0.29947492304906753,0.2028776978417266,0.4504038772213244,0.29045643153526973
18 0.28915444504798116,0.10935251798561152,0.43715670436187376,0.28215767634854777
19 0.30363932645301467,0.11223021582733814,0.4478190630048467,0.29460580912863077
20 0.28879232301285535,0.0,0.4500807754442649,0.28215767634854777
21 0.3306174180698896,0.0,0.4084006462035544,0.3360995850622407
22 0.5267065000905305,0.0,0.3066235864297255,0.5228215767634855
23 0.690747782002535,0.0,0.25137318255250396,0.6846473029045643
24 0.6291870360311426,0.0,0.2911147011308567,0.6307053941908715
25 0.4615245337678799,0.0,0.35508885298869153,0.4522821576763486
26 0.3981531776208583,0.0,0.37996768982229445,0.39834024896265563
27 0.2871627738547891,0.0,0.43489499192245556,0.27800829875518673
28 0.28607640774941157,0.0,0.4084006462035544,0.27800829875518673
29 0.2853521636791599,0.0,0.4006462035541195,0.27800829875518673
30 0.28806807894260367,0.0,0.39709208400646245,0.28215767634854777
31 0.28082563824008694,0.0,0.390306946688207,0.27385892116182575
```

補充

- 網頁顯示結果

Completed Applications (1)								
Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
<a href="#">app-20241016235426-0000</a>	Electric Power Consumption	16	1024.0 MiB		2024/10/16 23:54:26	Zhe	FINISHED	1.3 min