

Data Remediation by Spark for Web Crawler

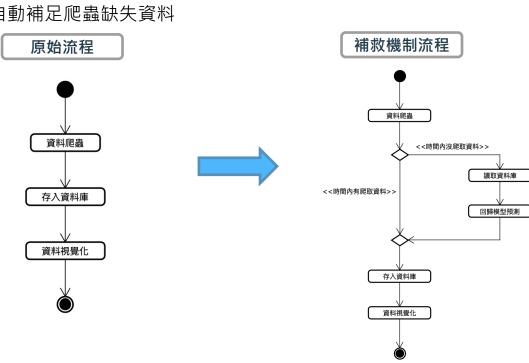


資碩計— 110753157 吳仁凱 資碩計一 張修誠 110753165 資碩計— 何彥南 110753202 莊崴宇 資碩工一 110753117 姚惠馨 資碩工一 110753135 資管四甲 107306009 吳泓澈

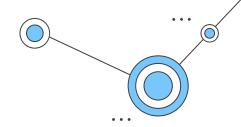
專案介紹

▶ 功能性:網路爬蟲之缺失值補救機制與視覺化

▶ 非功能性:自動補足爬蟲缺失資料











Web Crawler 爬取當日台積電股價資料



Data Prediction 回歸模型預測爬蟲遺失資料



Database 儲存爬蟲與預測資料

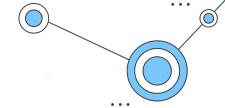


Data Visualization 資料視覺化

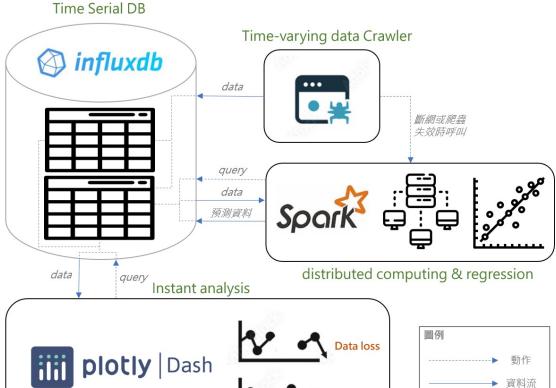


Demo Demo影片

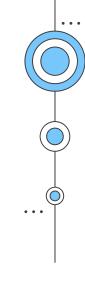
API展示圖



關聯







01 Web Crawler





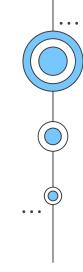
Web Crawler

- FinMind API
 - → 金融開源資料庫
- ▶ 時間序列資料
 - → 每二十秒call一次API
- ▶ 每日股票資訊
 - → 欄位資訊:

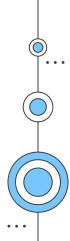
amount `average_price `buy_price `buy_volume `change_price `Change_rate `close `high `low `open `sell_price `sell_volume `total_amount `total_volume `volume `volume_ratio `yesterday_volume `date `stock_id `TickType

→ 以資料時間和收盤價進行視覺化





O2Data Prediction





MLlib

EX. 我們有一 個 1~12 的時 間序列輸入

His2 His1 Target

4

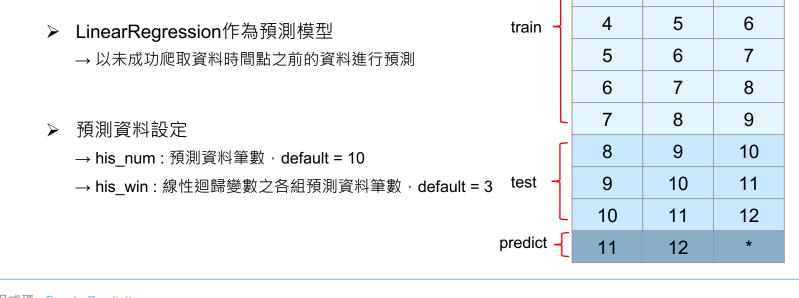
3

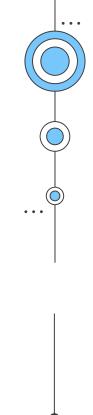
5

當爬蟲出現問題時

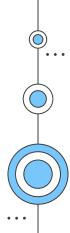
→ Call MLlib並傳送未成功爬取資料之時間







03 Database





InfluxDB

- ▶ 時間序列資料庫
 - → 適合儲存即時、不重複資料
- ➤ 爬蟲資料寫入兩個Table
 - → web_crawler_data存放爬蟲資料
 - → prediction_data存放爬蟲資料和預測資料
- > 支援多種程式語言
 - → 透過API即可直接取得資料庫資料

Read.py 說明

```
#引入influxdb套件
from influxdb import InfluxDBClient

#與DB建立連線
#InfluxDBClient(資料庫IP,資料庫PORT,娛號,密碼,DB名稱)
client = InfluxDBClient()

#捞取爬費TABLE的資料

prediction_data = client.query('select from web_crawler_data')

#捞取使費TABLE的資料

prediction_data = client.query('select * from prediction_data')

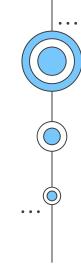
#對從爬費TABLE爬取的資料進行整理 並印出
print(list(web_crawler_data.get_points()))
print("-------")

#對從預測TABLE爬取的資料進行整理 並印出
print(list(prediction_data.get_points()))
print("-------")

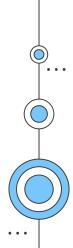
#對從預測TABLE爬取的資料進行整理 並印出
print(list(prediction_data.get_points()))
print(list(prediction_data.get_points())]

#對從預測TABLE爬取的資料進行整理(抓出最後一個層位的value) 並印出
print(list(prediction_data.get_points())[-1]['value'])
```





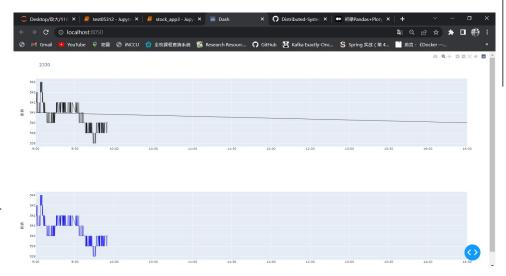
04 Data Visualization



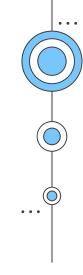


Visualization

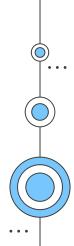
- Plotly API
 - → 將DataFrame繪製輸出視覺化圖表
- Dash API
 - → 將視覺化圖表部署至網站上
- ▶ 每五秒從Database更新一次資料







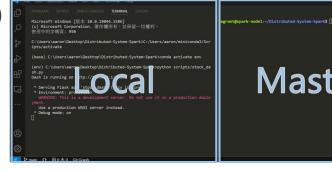
05 Demo





Demo Video





Master





Thanks!

Demo影片網址: https://www.youtube.com/watch?v=76fHSBQzw_Y

Github網址: <u>https://github.com/C-WeiYu/Distributed-System-Spark</u>

