

## 第二次讀書會報告

組別：11

組員：

邱昱榕 N114320004

朱怡樺 N114320005

林威呈 N114320011

董力慈 N114320019

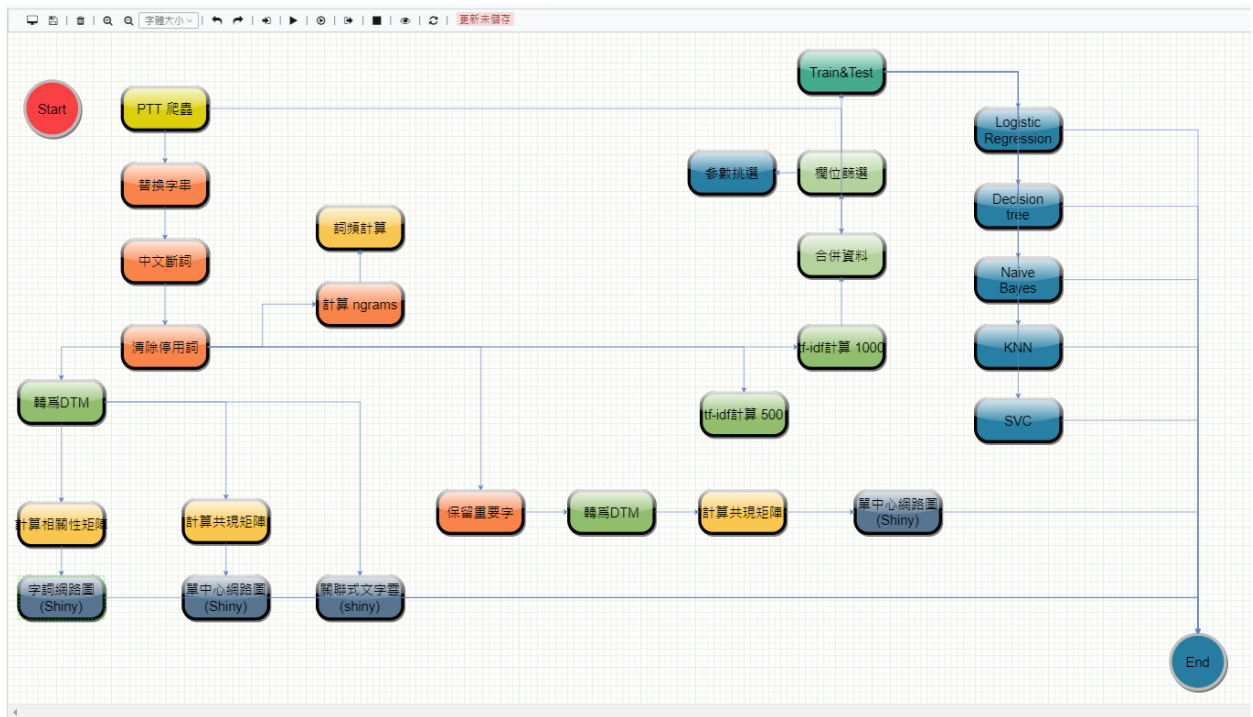
花秀旻 N114320021

趙承德 N114320023

陳昱維 N114320026

李晉安 N114320028

下圖為我們這次整個分類器訓練的流程圖



本次報告主題標的：

取 Tarflow 中，soft\_Job (軟體工作)、Stock(股票) & Tech\_Job ( 科技工作)

標的資料期間：

2024/01/01-2024/01/31

文字探勘工作流程設計平台

登錄test6

查詢與回報

≡

PTT 爬蟲 ( 27 )

✕

參數設定

任務結果

選擇看板 \*

Soft\_Job(軟體工作)

Steam(Steam)

Stock(股票)

studyteacher(實習教師)

TaichungBun(台中)

Tainan(台南)

TaiwanDrama(台劇)

搜尋關鍵字 ❶

以換行區隔，e.g.

國立中山大學

西子灣

...

排除關鍵字 ❶

以換行區隔，e.g.

壽山動物園

猴子

...

搜尋起始日期

2024/01/01

📅

搜尋結束日期

2024/01/31

📅

儲存更改

文字探勘工作流程設計平台

登錄test6

查詢與回報

≡

PTT 爬蟲 ( 27 )

✕

參數設定

任務結果

選擇看板 \*

Stock(股票)

studyteacher(實習教師)

TaichungBun(台中)

Tainan(台南)

TaiwanDrama(台劇)

Teacher(教師)

Tech\_Job(科技工作)

搜尋關鍵字 ❶

以換行區隔，e.g.

國立中山大學

西子灣

...

排除關鍵字 ❶

以換行區隔，e.g.

壽山動物園

猴子

...

搜尋起始日期

2024/01/01

📅

搜尋結束日期

2024/01/31

📅

儲存更改

(36) 針對 result 欄位做<中文斷詞>

定義詞彙的部分有<現金股利、股票股利、台積電...等等>

文字探勘工作流程設計平台 登錄test6 查詢問題回報

≡ 中文斷詞 (36)

參數設定 Input - 30 任務結果

選擇處理欄位 \*

result

定義詞彙 ⓘ

現金股利 1000  
股票股利 1000  
台積電 500  
聯發科 500  
瑞豐 500

選取字典 ⓘ

請選擇

儲存更改

(40) 使用功能<清除停用詞>，詳細設定可以參閱下圖，我們透過多次情緒分析結果的文字雲，來回觀察並增訂所需停用之字詞

選擇清除英文字母，因文章內容常用有網址、無效英文縮寫等

選擇清除數字，例如股票代碼等與情緒分析較無關之資料，以及技術性專有名詞等無情緒表達的字詞。

1. 常有文本的固定格式字詞，例如有依版規、內文、標題、資料來源等詞彙
2. 感謝詞包含謝謝、感謝等非表正面情緒的禮貌性用詞
3. 文章內表達時間序，例如過去、最近、現在、目前、未來、今天、今日、今年、去年、明年等詞彙
4. 常出現但與情緒分析較無相關之詞彙，例如成交量、股名、股票代碼、政府、產業等詞彙

清除停用詞 (40)

參數設定

Input - 36

任務結果

語言 \*

Chinese

使用預設停止詞

是

是否清除單字元 ⓘ

是

是否轉為小寫英文

是

清除英文字母 \*

是

清除數字 \*

是

清除換行符號 \*

是

清除特殊標點符號 \*

是

清除html tag \*

是

自定義停止詞

imgur  
JPG  
HTTP  
資料來源  
內文

儲存更改

分詞：將文本分割成單詞或字符的序列。

統計頻次：計算每個 n-gram 在文本中出現的頻次或概率。

最後輸出結果 ngrams 數量為 : 237848

通過詞頻計算，我們了解文本中各個詞彙的重要性和出現頻率，從而進行詞彙的分析和挖掘。此時我們利用詞頻計算所產生出來的文字雲，可以馬上了解本文章的重點

文字探勘工作流程設計平台

主 詞頻計算 ( 103 )

參數設定 Input - 102 任務結果

文字雲

根據詞彙表和詞頻信息，生成文檔詞頻矩陣（DTM），其中每一行代表一個文檔，每一列代表一個詞彙，單元格中的值表示對應詞彙在該文檔中的詞頻或其他統計信息，詞彙數量的值為1000

文字探勘工作流程設計平台 登錄test6 問題回報

### 轉為DTM (110)

參數設定

Input - 40

任務結果

**保留詞彙**  
以換行符號區隔，e.g.  
國立中山大學  
西子灣  
壽山...

**最多篩選詞彙數量**  
1000

儲存更改

文字探勘工作流程設計平台 登錄test6 問題回報

### 轉為DTM (110)

參數設定

Input - 40

任務結果

**統計資訊**

1001  
字數

2021  
文章數

**任務結果**

Show 10 entries Search:

system_id	一 下	一 堆	一 年	一 次	一 波	一 直	一 種	一 金	一 點	三 大	三 年	三 星	上 升	上 半 年	上 市	上 市 櫃	上 櫃	上 海	上 漲	上 班	上 述	下 午	下 半 年
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2	1.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
7	1.0	0.0	0.0	0.0	0.0	2.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2.0	0.0	0.0	0.0
8	2.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
10	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Showing 1 to 10 of 100 entries Previous 1 2 3 4 5 ... 10 Next

全量基瀏覽

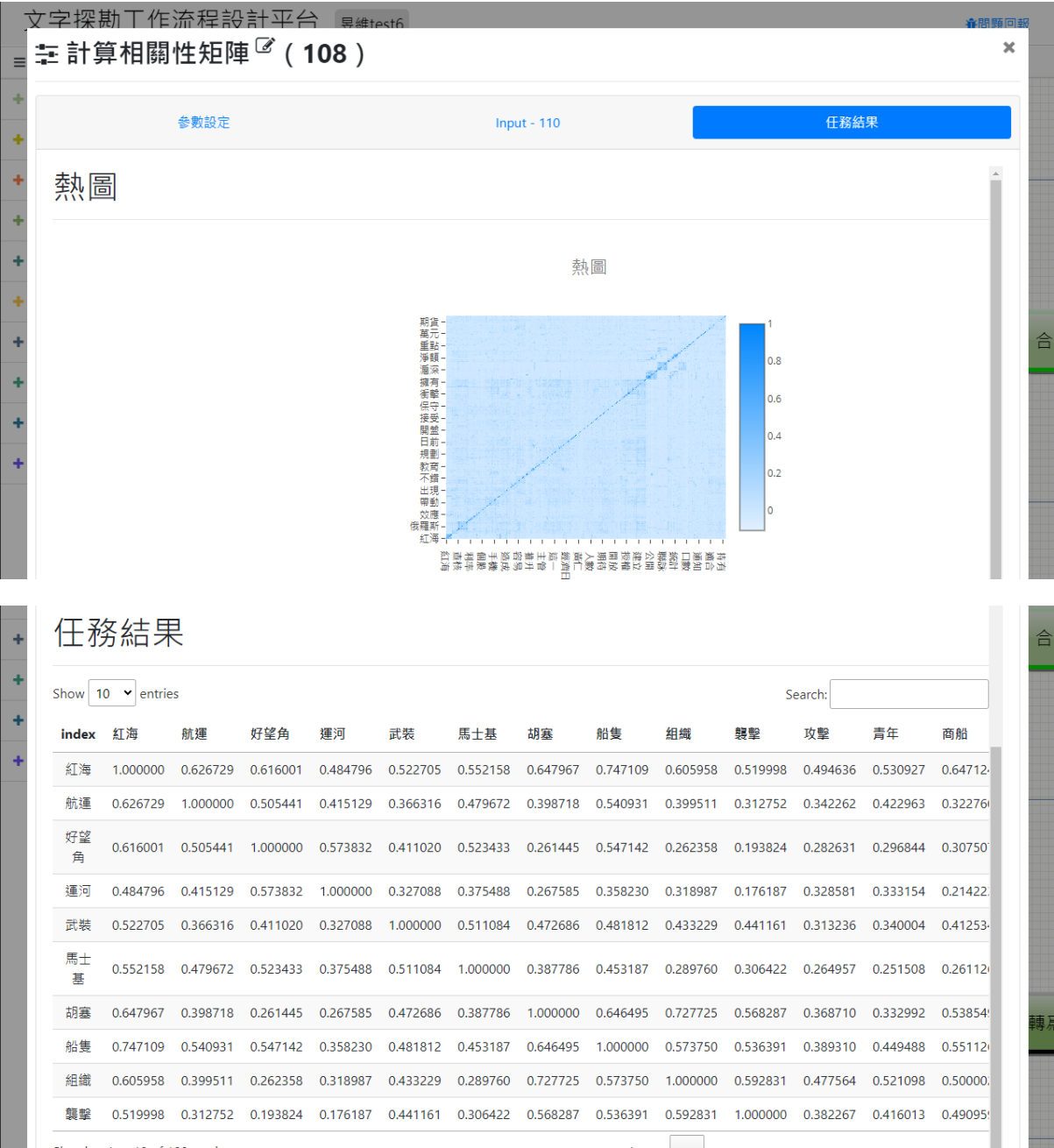
點我下載完整CSV資料

點我下載完整Rdata

點我下載完整json資料

計算相關性矩陣的目的是幫助我們理解文本數據中各個文本之間的相似度或關聯程度。這對於文本分類、主題建模、情感分析等任務非常有用，可以幫助識別相似主題、找到相似文本或文本群體等。

以下熱圖為展示數據之間的相似性、相關性和趨勢。



字詞網絡圖

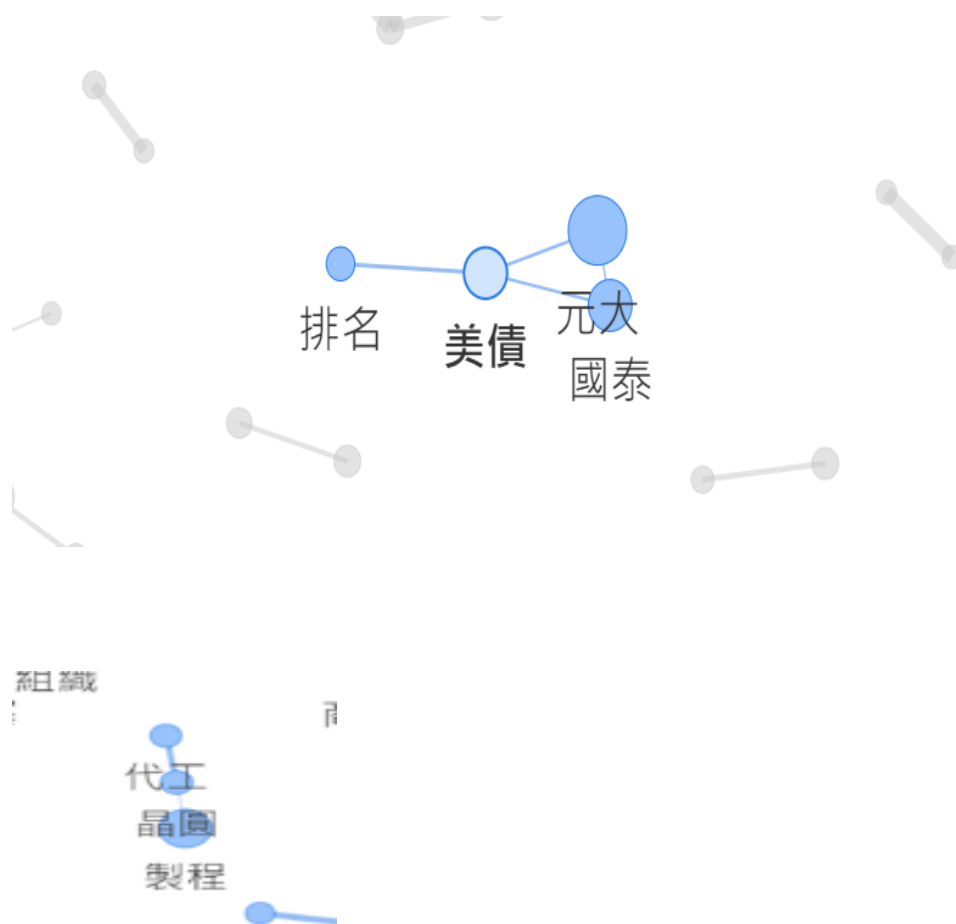
詞彙關聯距離

字體大小

本字彙圖所呈現的邊 (edges) 數量: 1  
本字彙圖所呈現的節點 (nodes) 數量: 80  
提示: 同一類別的字彙通常有相近的顏色

開始遊戲

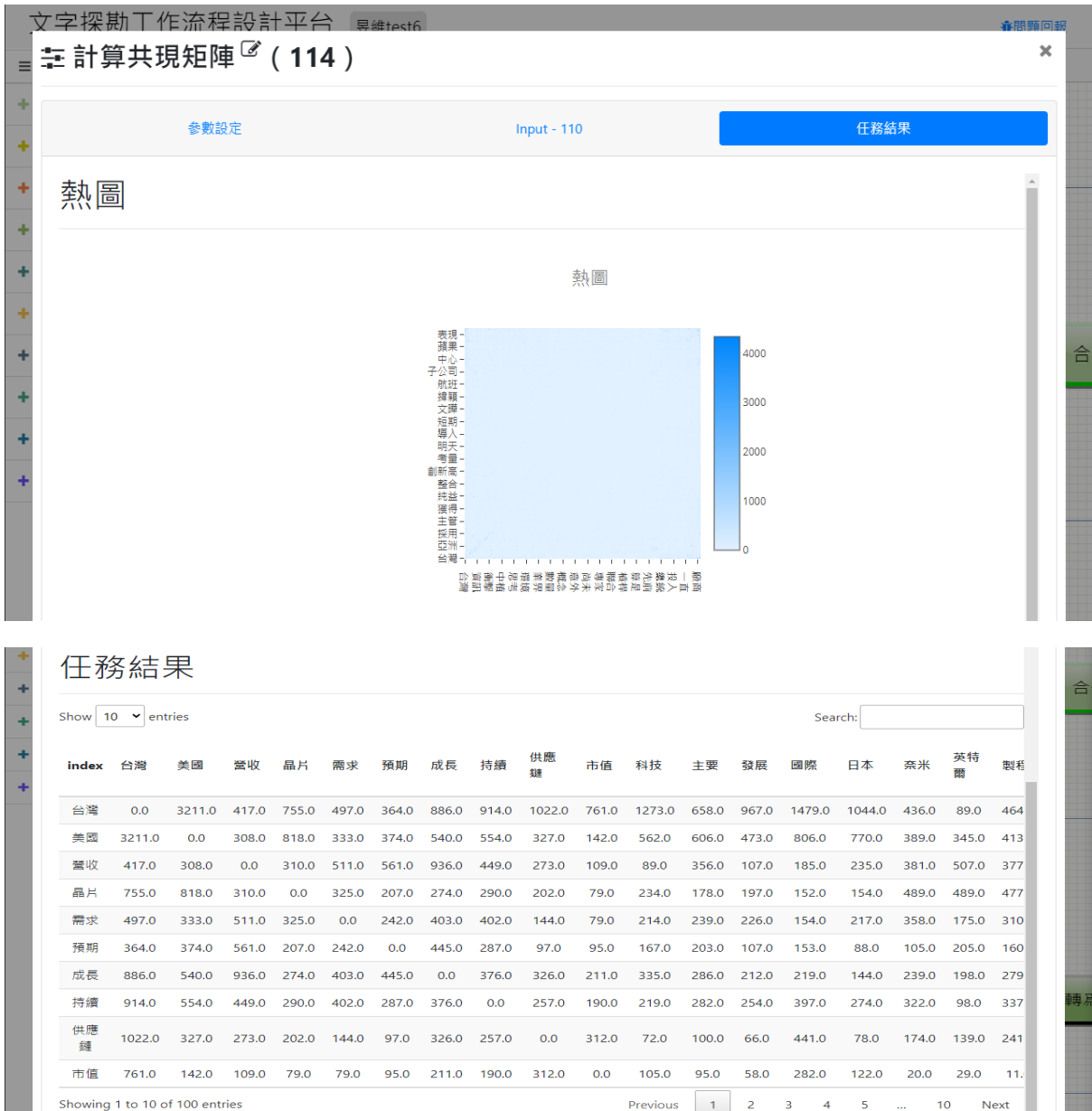
健身 課程 運算 量子 電腦 船隻 紅海 舉辦 旺年會 大學 教育 資產 管理 商船 葉門 取消 航班 地緣 政治 增資 新光金 緣由 發生 事由 相互 訊息 重大 採用 授權 購買 淨額 百分比 淨利 稅前 合併 觀測站 公開 平倉 口數 運動 青年 民眾 上半年 機電 廠區 綠電 網站 駭客 京鼎 製程 晶圓 代工 胡塞 組織 定期 定額 獎金 年終 人工 智慧 出口 管制 廠商 查核 合約 疫苗 採購 台指 台指期 收益 平準 櫃檯 百萬元 資本 公積 金控 機師 罷工 元大 美債 國泰 工具機 俄羅斯 融資 餘額 信用 券 增 統計



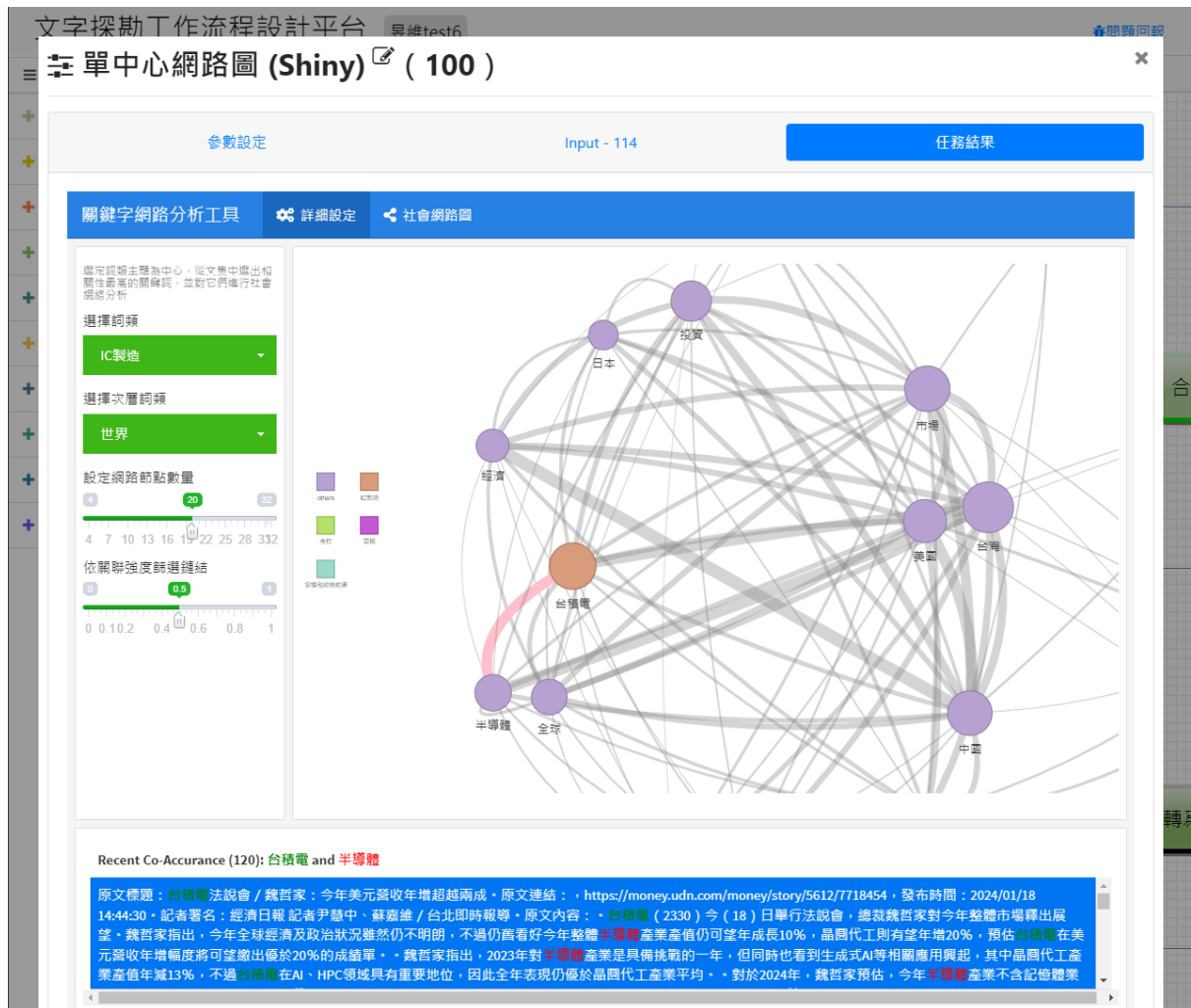




我們為了捕捉詞彙之間的語義關聯性,根據詞彙表和共現計算結果，生成共現矩陣，其中每一個元素表示對應詞彙之間的共現次數或其他統計信息。



我們先編輯了關鍵字點檔案，並選擇<IC 製造中>的<世界>來探索哪些字最常一起出現。其中<台積電>與半導體時常於文章中一起出現，也可於下方文章列表表中找尋。



我們要去把相同的字放在一起，並做一定的分類分群，下圖結果會將不同顏色代表不同分群，相比傳統文字雲，更傾向相同類別的字會放在一起

關聯式文字雲 (shiny) ( 121 )

參數設定 Input - 110 任務結果

分群數 \*

20

迭代次數(最少250次)

1000

聚合演算法 \*

weighted

文字雲顯示字數 \*

1000

距離計算公式 \*

euclidean



# 分類器

我們使用 **tf-idf**，用來評估一字詞對於一個文件集重要程度，字詞的重要性隨著它在文件中出現的次數成正比增加，但同時會隨著它在語料庫中出現的頻率成反比下降，我們選用 1000 個不同詞彙。

文字探勘工作流程設計平台 登錄test6 查詢問題回報

### tf-idf計算 1000 ( 134 )

參數設定

Input - 40

任務結果

**保留詞彙**  
以換行符號區隔，e.g.  
國立中山大學  
西子灣  
壽山...

**最多篩選詞彙數量**  
1000

儲存更改

文字探勘工作流程設計平台 登錄test6 查詢問題回報

### tf-idf計算 1000 ( 134 )

參數設定

Input - 40

任務結果

**統計資訊**

**1001**  
字數

**2021**  
文章數

**任務結果**

Show 10 entries Search:

system_id	一下	一堆	一年	一次	一波	一直	一種	一金	一點	三大	三年	三星	上升	上半年	上升
1	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
2	0.077028	0.000000	0.000000	0.000000	0.000000	0.079688	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
3	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
4	0.108801	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
5	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
6	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
7	0.097354	0.000000	0.000000	0.000000	0.000000	0.201433	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
8	0.454438	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
9	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.106701	0.000000	0.0	0.000000	0.000000
10	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.058810	0.000000	0.0	0.000000	0.000000

我們將原始資料與 TF-idf 跑出之特徵結果進行結合

文字探勘工作流程設計平台 登錄test6 查詢問題回報

合併資料 (137)

參數設定 Input - 27 Input - 134 任務結果

JOIN規則

新增規則 刪除規則

任務一欄位 system\_id 任務二欄位 system\_id

-----請選擇-----

並利用欄位篩選僅保留 System\_id、ArtCcatogory、其他的保留字詞與 TF-idf 跑出之特徵

文字探勘工作流程設計平台 登錄test6 查詢問題回報

欄位篩選 (138)

參數設定 Input - 137 任務結果

選擇要保留的欄位(按住ctrl(Windows)或command(MAC)可以複選) \*

system\_id  
artUrl  
artTitle  
artDate  
artPoster  
artCategory  
artContent

儲存更改

文字探勘工作流程設計平台 登錄test6 查詢問題回報

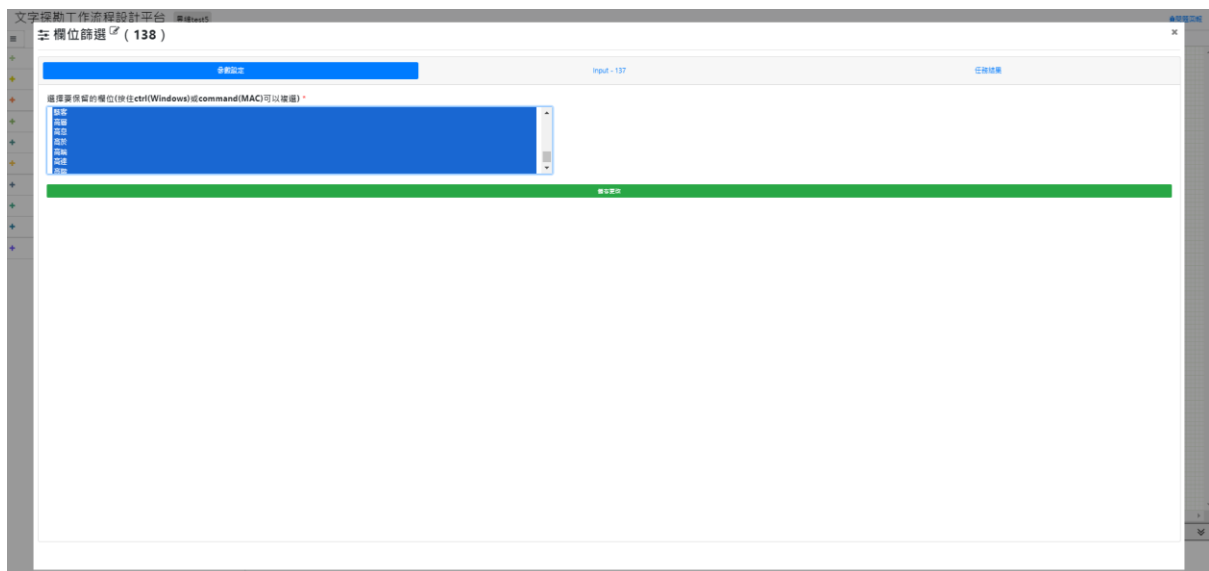
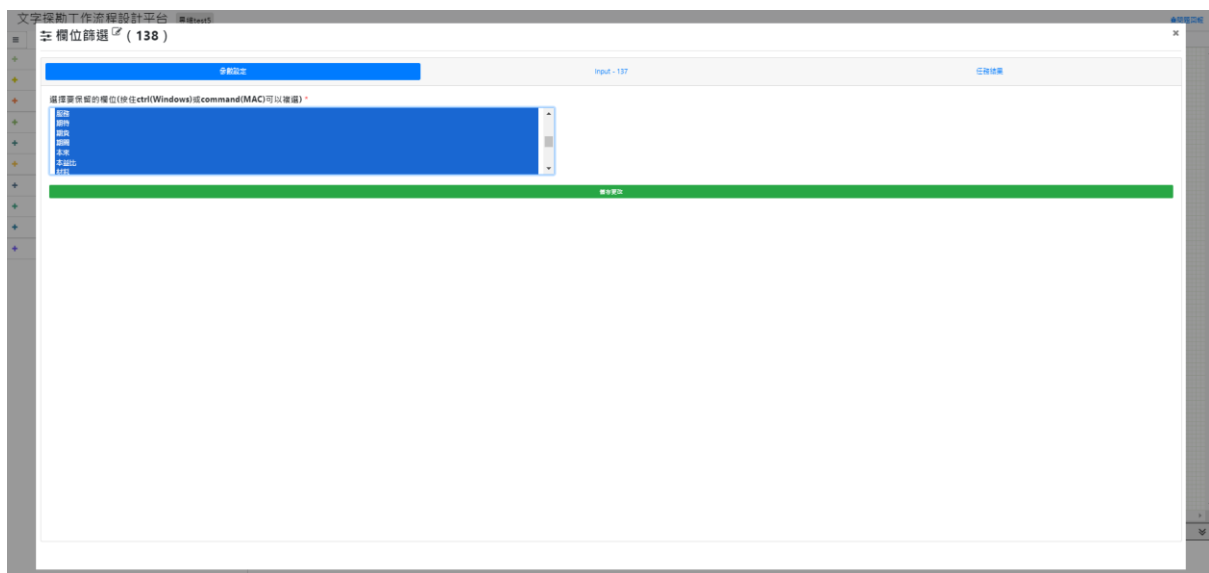
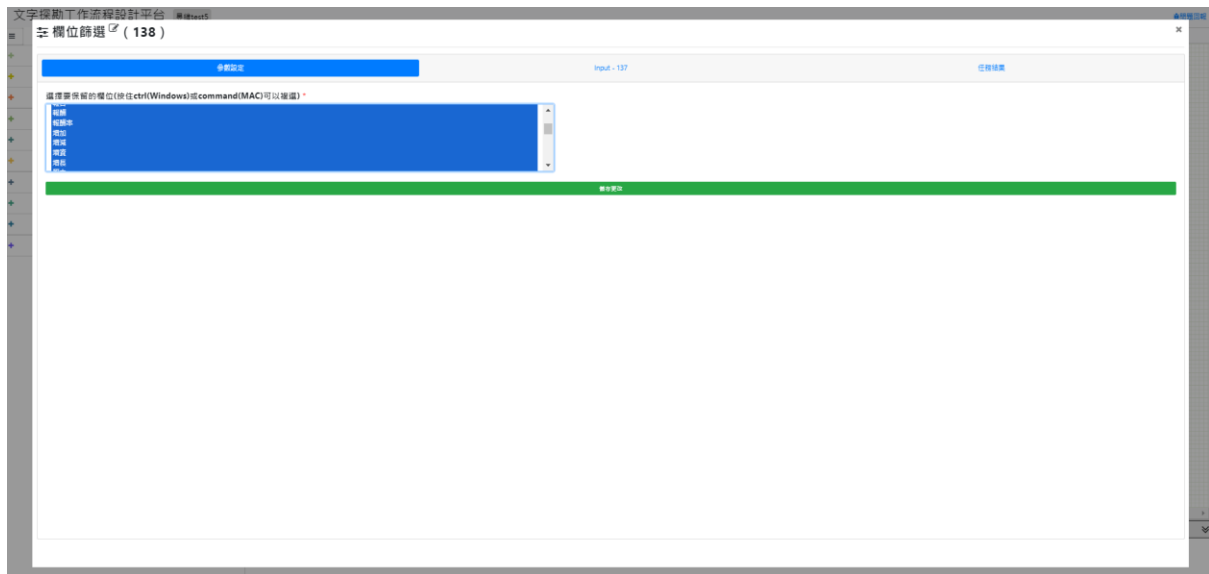
欄位篩選 (138)

參數設定 Input - 137 任務結果

選擇要保留的欄位(按住ctrl(Windows)或command(MAC)可以複選) \*

system\_id  
artUrl  
artTitle  
artDate  
artPoster  
artCategory  
artContent

儲存更改



我們依目標欄位去選擇切割比例，因原始資料是依序撈取，故將資料打散，並隨機抽取 20%的資料

文字探勘工作流程設計平台 呈維test6

### Train&Test ( 140 )

參數設定 Input - 138 任務結果

目標欄位 \*  
artCategory

測試資料切割比率 \*  
0.2

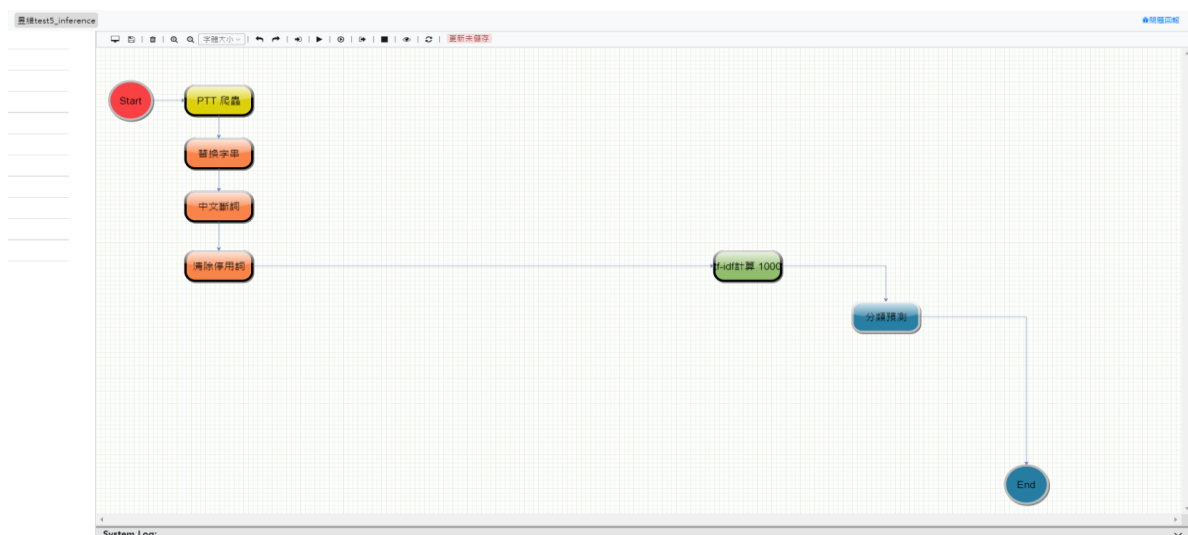
是否隨機排序資料  
是

亂數種子  
777

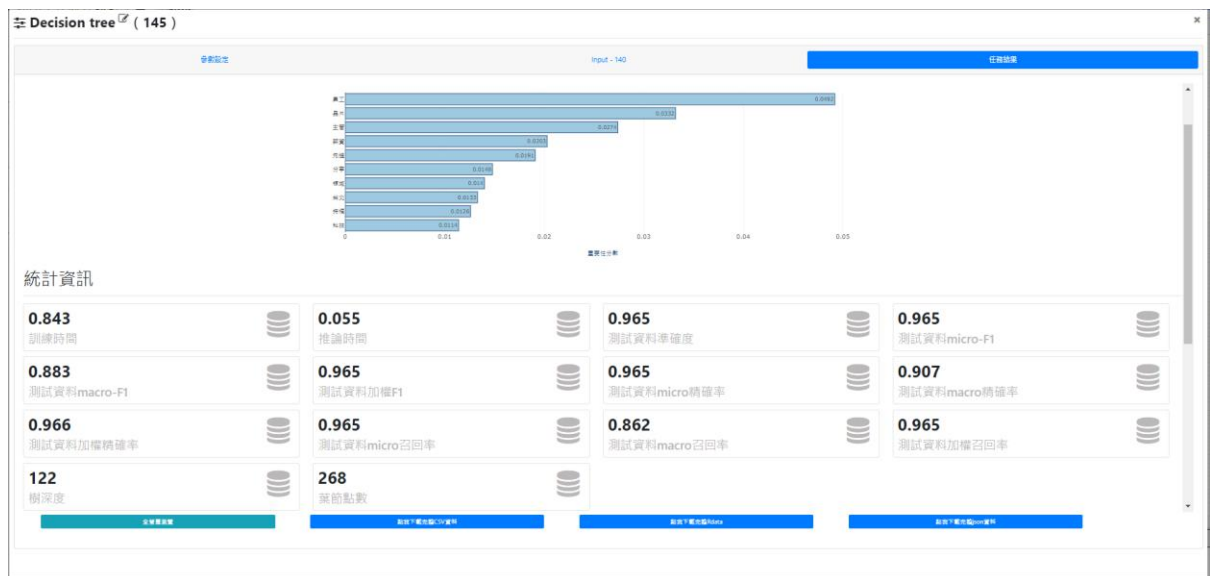
儲存更改



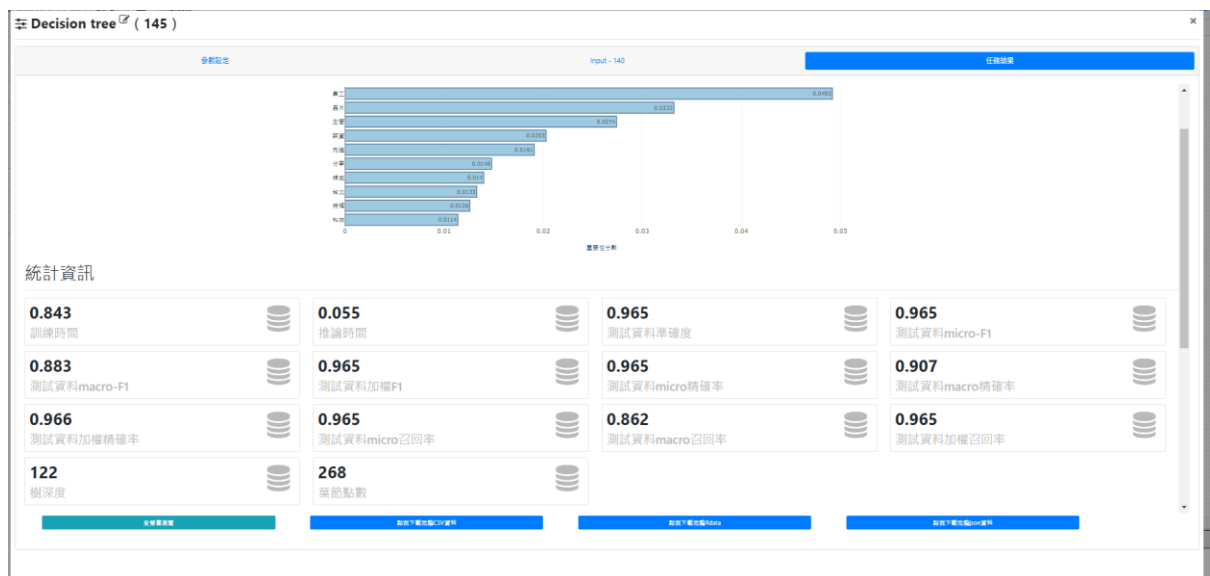
在 Logistic Regression 任務結果中，產出 inference workflow 的分類預測圖表功能







## 決策樹



## Navie Bayes



## KNN

字樣顯示工作並報設計平台

字樣顯示

主 KNN ( 147 )

參數設定

Input - 140

任務結果

統計資訊

0.033

訓練時間

0.49

測試資料macro-F1

0.831

測試資料加權精確率

5

鄰居數

0.346

推論時間

0.825

測試資料加權F1

0.847

測試資料micro召回率

euclidean

距離公式

0.847

測試資料準確度

0.847

測試資料micro精確率

0.468

測試資料macro召回率

0.847

測試資料micro-F1

0.551

測試資料macro精確率

0.847

測試資料加權召回率

## SVC

文字編輯工作系統設計平台 - 測試結果

≡ SVC (148)

參數設定

Input - 140

任務結果

統計資訊

1.689

訓練時間

0.811

測試資料macro-F1

0.943

測試資料加權精確率

1

懲罰係數

0.765

推論時間

0.941

測試資料加權F1

0.943

測試資料micro召回率

rbf

核函數

0.943

測試資料準確度

0.943

測試資料micro精確率

0.753

測試資料macro召回率

3

維度

0.943

測試資料micro-F1

0.945

測試資料macro精確率

0.943

測試資料加權召回率

參數欄位

Show 10 entries

Search