

國立東華大學應用數學系

111 學年度第二學期專題

指導教授：曹振海 博士

## 紅樓夢人物分析

*The Analysis of Characters in Dream of the Red Chamber*



學生：許堯智 李世勛 黃定綸 撰

中華民國一一二年六月

# 目錄

# 第一章 緒論

## 1.1 研究背景與動機

中國古典小說是中國文學中不可或缺的重要部份，在歷史上有許多值得進行探勘的名著。其中《紅樓夢》<sup>[1]</sup>被譽為中國古典小說的代表作之一，並且具有很高的探究性，小說中人物的命運和許多事件間有著密切而複雜的關係，需要經過深入的研究才能理解其真正含義。

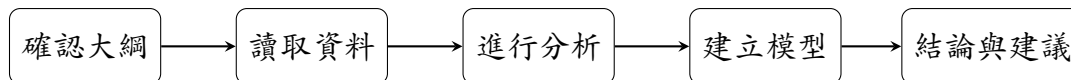
在探究過程中，我們參考了人物關係圖《紅樓夢|人物關係表|四大家族|寶玉、黛玉、寶釵血緣關係|5分鐘看懂紅樓夢人物關係|淺談紅樓夢|2023》<sup>[2]</sup>。在關係圖中可以見到人物關係錯綜複雜，最主要的四大家族即使有顏色標示，也難以明白其關係。但這已是研究者經過四天四夜建構出的關係圖，研究者也承認很難將人物關係完全弄懂。

雖然同樣存在著人物關係相對單純的小說，但正因《紅樓夢》值得深掘的研究價值，在文學界也誕生了新名詞「紅學」。在本研究中，《紅樓夢》被選為分析對象有其特殊的原因。首先，《紅樓夢》是一部長篇白話小說，使用的是比較容易理解的白話文，相對於使用大量文言文的古典名著，更易於進行文字分析和處理。其次，《紅樓夢》具有固定的主角群，自開頭到結尾都以同一群人推進劇情，這種連貫性使得人物之間的關係更為明確。

隨著生活節奏的加快，伴隨著大量的外部訊息和刺激，人們的專注力正逐漸下降。因此，為了吸引讀者的注意力，懶人包儼然成為一種新型流行趨勢，不僅使得訊息的呈現變得簡潔易懂，同時還具有足夠的視覺吸引力，並以一種簡單、直觀的方式呈現訊息，貼近了現代人快節奏的生活步調。本研究旨在希望能透過分析《紅樓夢》中人物出現的次數、頻率，並以簡單、直觀的方式呈現人物間的關係和生活圈，使讀者能夠藉由本研究了解《紅樓夢》主要角色的關係。作為一部白話章回小說，《紅樓夢》提供了明確的人物關係和通俗的文字敘述，為研究和分析提供了有價值的案例。本研究除了以能讓讀者更深入地了解小說的人物關係為目的，同時也以滿足現代社會對於簡潔、易懂和吸引人的訊息呈現的需求作為目標。

## 1.2 研究流程

本次研究流程與架構分爲以下五步驟：



圖一 研究流程架構圖

研究流程說明：

### 確認大綱：

從本身所接觸到的各種領域及應用中，選擇文字探勘技術作為分析工具，並藉由文獻探討現有的相關研究方向，進一步決定研究主題。最後，將數據資料以清晰、直觀的方式呈現。

### 讀取資料：

使用 R 語言將《紅樓夢》文本轉換成可進行分析的格式，例如進行分詞、刪除標點符號，並將人物及其別名等資料抓取出來，以及對資料做初步的處理。

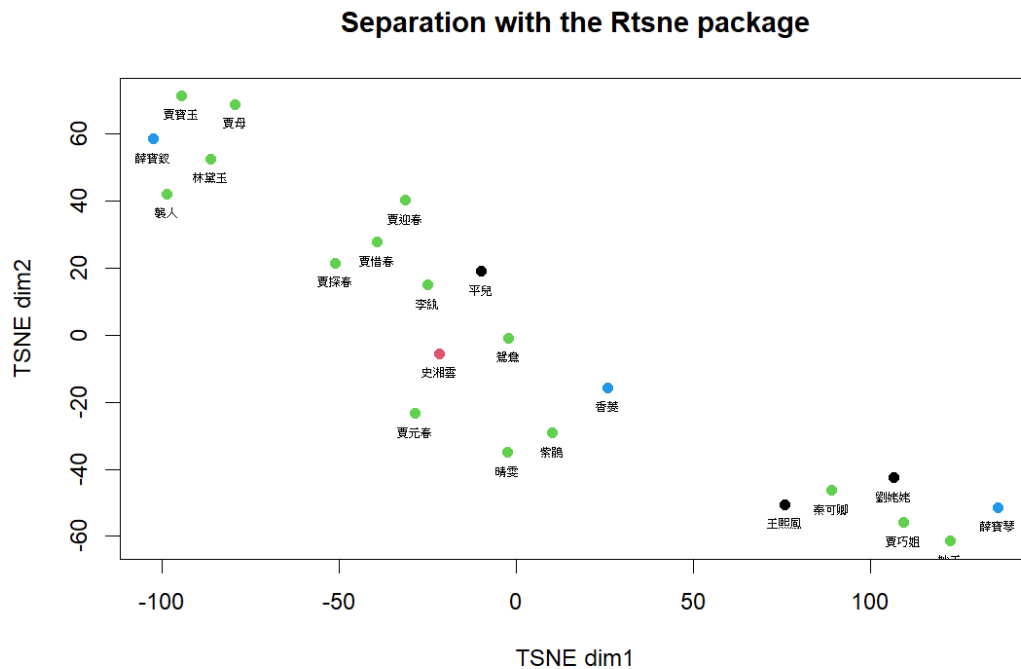
<MINTED>

### 進行分析：

透過對資料的分析繪製可視化圖表，如次數長條圖、相關係數矩陣以及角色之間的相關性，並針對資料與圖表進行初步分析，以提供作為後續計算模型依據。

### 建立模型：

利用 networkD3、Rtsne 套件，最終做出的人物關係進行可互動的視覺化網路圖，可將研究成果放在網路上供比對與參考。



圖二 t-SNE 分群

### 結論與建議：

不同的模型具有不同的可解釋性，本研究會將結果與原先預測來做比對，並加以整理之後作為本研究的結論。不僅可以獲得分類結果，還可以獲取人物關係圖等重要訊息，從而更進一步的分析和解釋，獲得更多有用的資訊。

除此之外，也會將研究過程中本應還可再進行深入的部份，一同建議研究方法，希望能讓下一位研究者獲得靈感。

## 第二章 文獻研究

### 2.1 分析方式簡介

以下為本研究中使用到的分析方法：

#### **correlation：**

相關係數使用數值表示，是用來描述兩個變數之間相關程度的指標，表示兩個變數相互關聯的程度，其取值範圍介於 -1 到 +1 之間，其中正的相關係數表示兩者呈現正相關，負的相關係數表示兩者呈現負相關，而接近零的相關係數表示兩者之間沒有關係。相關係數可以幫助我們了解兩個變數之間的關聯性，並發現資料數據的發展趨勢。

#### **wordcloud：**

文字雲是一種用於文字資料視覺化的方法。這種方法可以將大量文字資料轉換成一個具有可視性的圖形，將高出現頻率的詞彙以較大的字型呈現，並集中在中央；低出現頻率的詞彙以較小的字型呈現，會散佈於外圍，是資料的一維呈現方法。

#### **MDS：**

MDS 全名為 Multi-dimensional Scaling，是一種可將高維度的資料降至二維或三維平面上，以更易懂地理解和分析數據的方式。通常用於探索數據之間的相似性或差異性，並以圖形化的方式展示數據的結構。在《紅樓夢》人物關係的研究中，MDS 可把人物之間的相似度轉換為距離單位，然後將這些距離單位映射到二維或三維空間中，並將每個人物表示為一個點，人物之間的距離表示他們之間的相似程度。

#### **Jaccard：**

Jaccard 用於評估兩個文本之間的相似程度，這可以通過將兩個文本中共同出現的單詞數量除以兩個文本中所有不同單詞的總數來完成。簡單來說分子是交集，分母則是聯集。在 R 語言中使用該函數，可以更精確計算兩個文本之間的相似度，進而進行文本分析和比較。

### PCA：

PCA 把高維度的資料以盡可能保留其特性的方式降到低維度，爲了將一組相關變量轉換爲一組新的不相關變量，稱爲主成分 (PC)。主成分是原始變量的線性組合併且彼此正交。第一個 PC 捕獲數據中的最大變異數，每個後續 PC 捕獲最大剩餘變異數，但前提是它與之前的 PC 正交。其降維方法主要利用線性的方式，也就是降維生成的每一個 PC 都是原本變數 ( $x_1, x_2, x_3 \dots$ ) 的線性組合。

### t-SNE：

t-SNE 將點之間的相似度轉化爲條件機率，原始空間點的相似度由常態分佈表示，嵌入空間中點的相似度由 t 分佈表示。t-SNE 的主要優勢在於解決了降維後的擁擠問題，使得相似的樣本能夠聚集在一起，而差異大的樣本能夠有效地分開，避免了其他降維方法各個點分佈擁擠、邊界不明顯的缺點

## 2.2 論文與相關研究參考

爲了尋找適合的視覺化結果，我們在碩博士論文網站中找到了《金庸小說互動式視覺化文字探勘》論文<sup>[?]</sup> ，論文提供的方法是文字雲：取出需要的人物名稱並合併別名後，可發現到當人物出現次數越多，則該人物字體越大，鮮豔又淺顯易懂，符合作爲懶人包的必要條件。只可惜文字雲侷限在單一人物的出現頻率，無法知道一對一或多對一之間的人物關聯。

接著，在知乎網站發掘到能顯示人物關聯性的網路圖<sup>[?]</sup> 。參考該網站後，發現該圖以人物爲節點，關聯性作爲線段粗細的網路圖，若各節點線段愈粗，便代表兩節點關係愈深，節點大小代表的是人物出現的頻率，顏色則代表的是人物所代表的家族，比起文字雲多了角色跟角色之間的關聯。

即使網路圖解決了人物與人物的關聯性問題，依舊還侷限在單一人物的對應關係，難以明確知道多人對應多人的關係，因此我們透過其他視覺化方式，來分析呈現資料特性。

爲此，從人物關係的網路圖轉而進行分群方法：MDS、PCA 跟 t-SNE。其中，分群之後的結果仍然較爲鬆散，爲了改進這方面的問題，在原本的分群套件上再加上名爲 Jaccard<sup>[?]</sup> 的篩選方式。

## 第三章 研究方法

### 3.1 使用套件

`corrplot`：描述不同變數之間相互關聯性的關係矩陣。

`tm`：在 R 中進行文本分析、處理、建模等任務變得更加方便，支持生成文字雲、頻率圖、主題分布圖等。

`wordcloud`：根據詞彙的頻率或權重，創建具有不同大小和顏色的文字雲圖形。

`jiebaR`：用於 R 語言中的中文分詞套件，將中文文本進行分詞處理，也是進行中文分詞最常見的方式。

`networkD3`：創建互動性網路圖可視化。

`readxl`：用於讀取和解析 Excel 文件。

`tidyverse`：包括 `ggplot2`、`dplyr`、`tidyr` 等等，涵蓋了分組統計、數據可視化、轉換等各方面的套件。

`igraph`：網絡圖的操作和轉換，這次用於創建邊數據。

`magrittr`：使得連續的函數調用和數據處理更加易讀和易寫，例如 `%>%`。

`dplyr`：用於對資料進行快速、直觀和一致的操作。

`ggpubr`：在 `ggplot2` 的基礎上構建的擴展套件，可視化功能和統計分析工具。

`MASS`：提供了線性模型、分類分析等的數據分析和建模。

`vegan`：用於處理和分析生態學數據，這次是使用 Jaccard 進行修正。

Rtsne：將高維數據降低到二維或三維，以便更好地理解數據的結構和關係。

## 3.2 文字探勘

文字探勘是從大量文字中提取感興趣資訊或挖掘有用知識，並透過電腦的運算能力，過濾及轉化大量的文字內容，找出隱含且有用的資訊。此處結合語言處理、統計分析等方法，於文本中找尋趨勢和關聯性。

由於《紅樓夢》擁有許多個性鮮明的角色，在此使用 R 語言將主要的角色與其別名進行整理，找出他們在文本中的出現頻率，以便後續分析不同角色間的互動模式，提高後續資料分析和模型建立的準確性和可靠性。

<MINTED>

`name_matrix` 為串聯所有章節中出現人物的名稱，其作用為方便做文字探勘，以方便後續資料處理。

<MINTED>

給定一空矩陣，透過 `stat_name` 這一個自定義的文文字探勘函數計算各個人物在各回的出現次數，並將人物名稱與出現次數合併到 `name` 矩陣。

<MINTED>

以下程式碼為 `stat_name` 函數的運作方式，此處作法部分參考自銘傳大學論文<sup>[?]</sup>，並加以修改成符合本研究所需的樣式。

<MINTED>

將 `name` 矩陣中屬於同一人物的名字與出現次數合併，統整為人物各個章回出現次數矩陣，以利進行後續的各項統計與分析。

<MINTED>

經由以上文字探勘與簡化的過程，我們得到以下人物總出現次數矩陣，此矩陣可運用在文字雲的呈現，並運用數據資料進行分析。

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
賈寶玉	4	2	32	0	61	14	29	85	35	5	10	1	11	16	57	25	48	41	117	56	49	43	54
林黛玉	0	1	83	3	10	0	9	25	4	0	0	3	2	1	0	8	0	18	33	24	14	29	18
薛寶釵	0	0	0	2	5	0	12	20	1	0	0	0	0	0	0	1	0	10	2	21	8	24	2
賈元春	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	4	5	44	3	2	0	1	3
賈探春	0	1	4	0	1	0	2	0	0	0	0	0	0	0	0	0	0	3	0	1	0	3	4
史湘雲	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	7	15	16	0
妙玉	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	3	0	0	0	0	0
賈迎春	0	1	3	1	1	0	2	0	0	0	0	0	0	1	0	2	0	2	0	2	0	5	2
賈惜春	0	1	1	0	1	0	8	0	0	0	0	0	0	0	0	0	0	1	0	1	0	2	3
王熙鳳	0	0	11	0	0	2	0	0	1	1	1	1	1	1	0	0	0	1	0	1	0	1	0
賈巧姐	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
李紈	0	0	0	3	0	0	4	0	0	0	0	0	0	0	0	1	0	5	0	0	1	3	0
甄可卿	0	0	0	0	16	1	13	2	0	8	18	0	14	1	1	0	0	0	0	0	0	0	0
賈母	0	0	33	3	6	2	3	15	7	0	6	4	4	0	3	10	3	27	3	4	3	42	7
劉姥姥	0	0	0	0	0	61	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
香菱	0	0	0	0	0	0	6	0	0	0	0	0	0	0	0	7	0	0	0	1	0	0	0
平兒	0	0	0	0	0	12	5	0	0	0	5	1	3	2	0	9	0	0	0	0	28	0	0
晴雯	0	0	0	0	1	0	0	6	1	0	0	0	0	0	0	0	0	0	1	6	0	0	0
襲人	0	0	9	0	4	11	0	6	5	0	0	0	2	0	0	0	0	1	55	14	24	9	7
紫鵲	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	3	0	0
麝香	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
薛寶琴	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

表一 name 矩陣/人物總出現次數矩陣

爲了能夠從資料裡提取到人物間的關係，此處我們將 name 矩陣再度簡化，僅表示該章回中人物是否有出現，並將結果儲存至 person\_exist 矩陣，稱之爲人物出現矩陣。

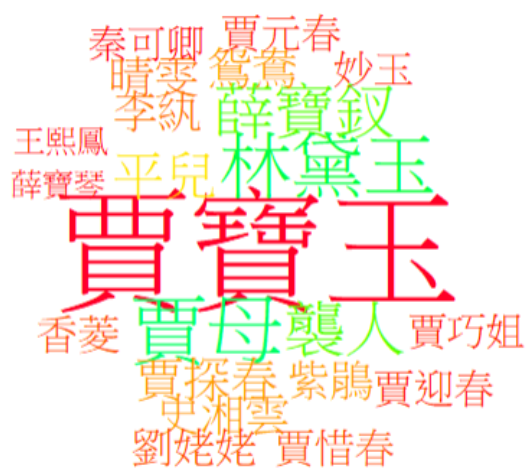
<MINTED>

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
賈寶玉	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
林黛玉	0	1	1	1	1	0	1	1	1	0	0	1	1	1	0	1	0	1	1	1	1	1	1
薛寶釵	0	0	0	1	1	0	1	1	1	0	0	0	0	0	0	1	0	1	1	1	1	1	1
賈元春	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	0	1	1
賈探春	0	1	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	1
史湘雲	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	0
妙玉	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
賈迎春	0	1	1	1	1	0	1	0	0	0	0	0	0	0	1	0	1	0	1	0	1	0	1
賈惜春	0	1	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	1
王熙鳳	0	0	1	0	0	1	0	0	1	1	1	1	1	1	0	0	0	1	0	1	0	1	0
賈巧姐	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
李紈	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	1	0	1	0	0	1	1	0
甄可卿	0	0	0	0	1	1	1	1	0	1	1	0	1	1	1	0	0	0	0	0	0	0	0
賈母	0	0	1	1	1	1	1	1	1	0	1	1	1	0	1	1	1	1	1	1	1	1	1
劉姥姥	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
香菱	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0
平兒	0	0	0	0	0	1	1	0	0	0	1	1	1	1	0	1	0	0	0	0	1	0	0
晴雯	0	0	0	0	1	0	0	1	1	0	0	0	0	0	0	0	0	0	1	1	0	0	0
襲人	0	0	1	0	1	1	0	1	1	0	0	0	1	0	0	0	0	1	1	1	1	1	1
紫鵲	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
麝香	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
薛寶琴	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

表二 人物出現矩陣

### 3.3 資料視覺化

資料視覺化是將數據以圖形或圖表的形式呈現，幫助人們將繁雜的數據簡化成易懂好吸收的內容。在 R 語言中，有多個套件可用於資料視覺化，包含 wordcloud 及 networkD3 等等。wordcloud 能讓讀者在不閱讀所有文章的前提下，快速了解並聚焦大批文章中主要的議題；而 networkD3 則可以透過顏色、節點大小或者線段粗細等等方式來展示資料的不同屬性，其可互動性也高於其他二維圖形。



圖三 使用 wordcloud 後的視覺化結果

## 第四章 研究探討與分析

### 4.1 數據清理

將文本讀入編譯器後，透過先前自製的篩選函數，蒐集《紅樓夢》各個篇章內容，計算每個人物在各章節總出現次數。我們可以篩選出各個人物的名稱，接著建立一個矩陣，用來紀錄每個人物是否出現在各回中，分別利用 1 跟 0 去判斷人物是否出現在該章回裡，這個矩陣將成為後續所有視覺化結果的資料來源。

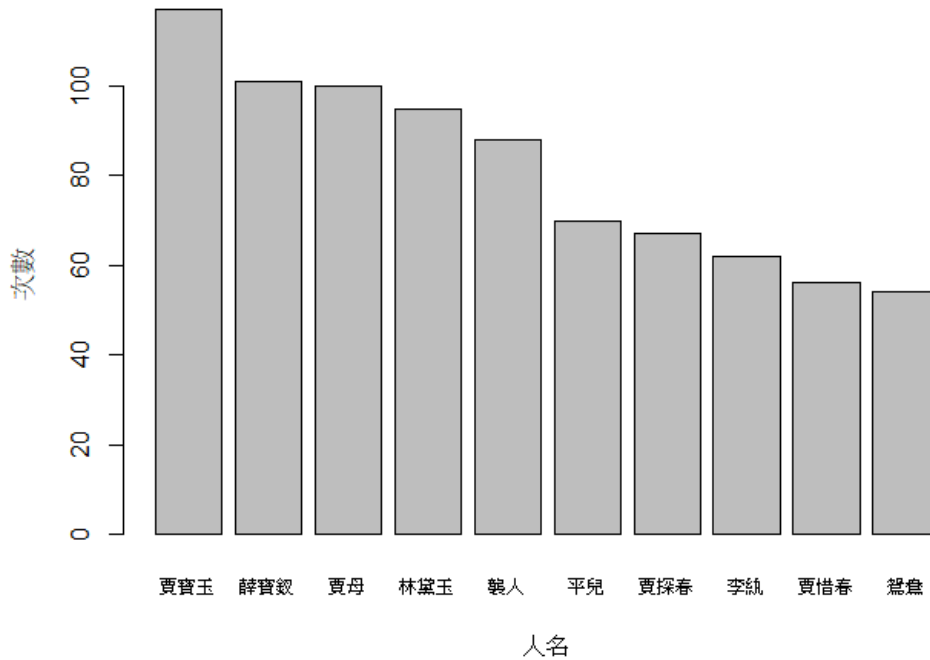
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
賈寶玉	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
林黛玉	0	1	1	1	1	0	1	1	1	0	0	1	1	1	0	1	0	1	1	1	1	1	1
薛寶釵	0	0	0	1	1	0	1	1	1	0	0	0	0	0	0	1	0	1	1	1	1	1	1
賈元春	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	0	1	1
賈探春	0	1	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	1
史湘雲	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	0
妙玉	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
賈迎春	0	1	1	1	1	0	1	0	0	0	0	0	0	1	0	1	0	1	0	1	0	1	1
賈惜春	0	1	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	1
王熙鳳	0	0	1	0	0	1	0	0	1	1	1	1	1	1	0	0	0	1	0	1	0	1	0
賈巧姐	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
李執	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	1	0	1	0	0	1	1	0
索可卿	0	0	0	0	1	1	1	1	0	1	1	0	1	1	1	0	0	0	0	0	0	0	0
賈母	0	0	1	1	1	1	1	1	1	0	1	1	1	0	1	1	1	1	1	1	1	1	1
劉姥姥	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
香菱	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0
平兒	0	0	0	0	0	1	1	0	0	0	1	1	1	1	0	1	0	0	0	0	1	0	0
晴雯	0	0	0	0	1	0	0	1	1	0	0	0	0	0	0	0	0	0	1	1	0	0	0
襲人	0	0	1	0	1	1	0	1	1	0	0	0	1	0	0	0	0	1	1	1	1	1	1
紫鵑	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
鴛鴦	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
薛寶琴	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

表三 人物出現矩陣

接著使用 R 語言中的相關套件來製作長條圖，根據出現次數從多到少進行排序，顯示各個人物包含別名在《紅樓夢》中的出現次數。賈寶玉、薛寶釵、賈母和林黛玉等人將在這張圖中佔據較前面的位置，表示他們是故事中重要的角色。

這張長條圖為我們提供一個直觀的方式，能更簡單地理解《紅樓夢》中各個人物的重要性和出現頻率。

第1回至第120回前10名人物出現次數統計表

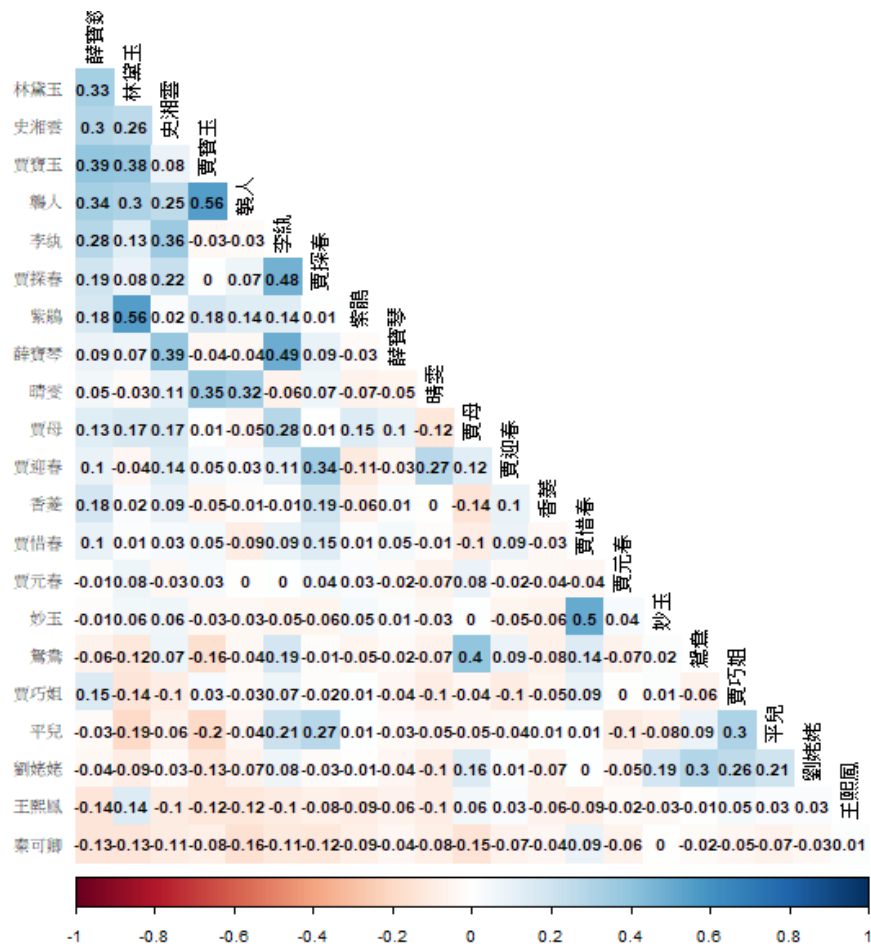


圖四 人物出現次數統計長條圖

利用上述人物出現矩陣進行初步視覺化，使用 corrplot 進行矩陣內所有角色的相關分析，其中角色之間的相關性將以數值和顏色的方式呈現。如果相關係數的絕對值愈接近 1，顏色會越深，表示相關性大；相反，相關係數愈接近 0，顏色則會越淺，表示相關性小。

以賈寶玉和林黛玉為例，他們之間的相關係數為 0.38，表示兩者之間相較於其他角色，存在密切的關係。而劉姥姥和王熙鳳的相關係數接近 0，表示他們與其他角色之間可能沒有太大關聯性。

此相關矩陣圖不僅能夠簡單分析各個人物之間的關聯性，還可以作為接下來 network、MDS、PCA 模型的對照。透過比較這些視覺化方法的結果，可以了解它們在展示和解釋數據方面的差異。

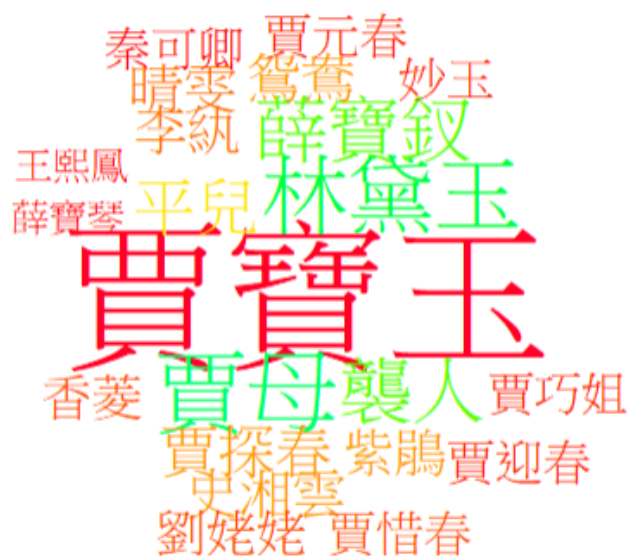


圖五 相關係數矩陣圖

## 4.2 文字雲

首先做的是人物總出現次數文字雲，除了參考上述論文<sup>[2]</sup>的文字雲模型，並以此為基礎進行修改，透過字體的大小來去呈現出此關鍵字的重要性，更可以一眼看出各個關鍵字的熱度。

其中人物出現次數愈多，則該人物字體愈大，且出現次數愈接近顏色也愈接近。但文字雲僅能透過字體的大小來去呈現出此關鍵字的重要性，無法藉此了解人物間彼此的關係，且當出現同樣文字大小的關鍵字，會因為關鍵字的字元長度不同，可能讓人產生誤判的狀況。



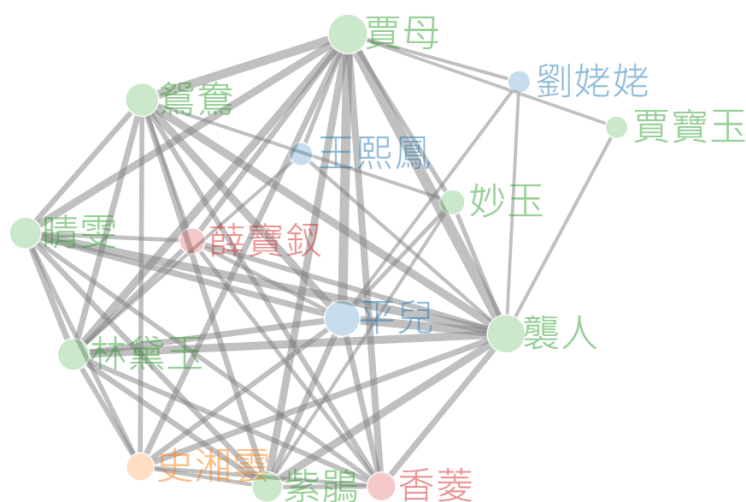
圖六 人物總出現次數文字雲

### 4.3 網路圖

由於文字雲難以顯示出人物關聯性的狀態，透過《R 語言基於共現提取《雪中悍刀行》人物關係並畫網絡圖》<sup>[7]</sup> 這篇網路文章，發現此程式碼生成的網路圖擁有高互動性，並且相對直觀的了解該人物之間的密切關係。以下將以此篇文章所介紹的程式碼進行繪製網路圖。

首先，將先載入主要的套件「networkd3」來繪製網路圖，便能將資料透過文章中提供的程式碼來呈現。

其中，該網站的程式碼是使用一種分詞器，一種名為「jiebaR」的套件來進行分詞。這個 R 語言的內建套件支援四種不同的分詞模式，同時還支援姓名、簡體中文、正體中文和關鍵字等功能。



圖七 網站程式碼所構出的網路圖

在此網路圖中，節點所代表的是該人物出現的頻率，節點間的線段代表人物在同一章回中共同出現的次數，其中線段愈粗者，表示人物間一起出現次數愈高，而節點顏色所代表的則是個人物所代表的家族。

然而，在網路文章中提供的程式碼執行後，發現結果和最初預期有相當大的區別，並且與人物出現次數統計長條圖、相關係數矩陣圖以及人物總出現次數文字雲這三張先前所製作出的圖形有資料上的落差。例如，該程式碼將「賈母」這個人物判定為《紅樓夢》中出現頻率最高的角色，但根據之前收集的資料和前面所繪製出的文字雲，出現頻率前三名依序應該是「賈寶玉」、「賈母」和「林黛玉」才對。

經過對程式碼的逐一排查，發現使用「jiebaR」套件的分詞器可能對於《紅樓夢》這種接近白話文但仍帶有部分文言文的作品分詞依然不夠精確，例如林黛玉的分詞「黛玉」可能會被變成諸如「和黛玉」、「比黛玉」、「黛玉同」等字詞，使得在繪製網路圖時，網路圖容易發生人物關係不正確的狀況。

由於上述所發生的問題，必須對分詞的資料進行了修正，於是以下使用之前建立的「人物出現矩陣」進行了篩選，並將兩個人物在同一回內出現的關係設置為1，否則設置為0。再根據這些關係計算頻率和相關資料重新繪製出如圖八的網路圖。

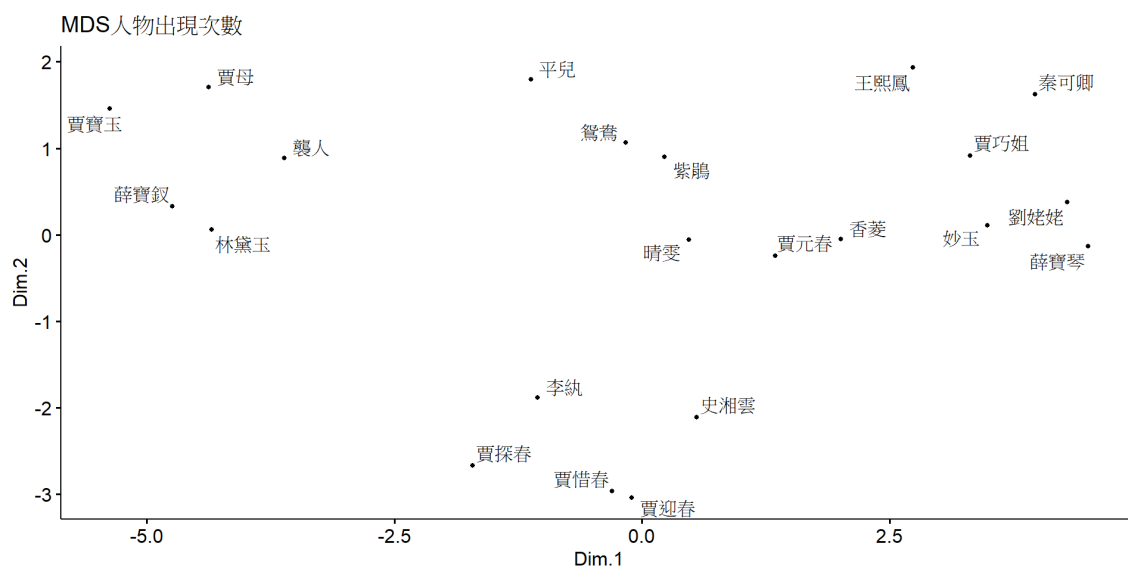


從上圖中，可以看到只繪製了主要互動的人物關係，並能輕鬆地透過節點大小與線段粗細分別看出人物的出現頻率和關係。然而，圖中人物之間的距離沒有具體的含義。爲了使人物之間的距離具有意義並更好地進行分群，於是採用了多元尺度分析（MDS）方法。

## 4.4 MDS

MDS 全名爲 Multi-dimensional Scaling，是一種多變量分析技術，也是一種降維方法，可以將高維度的資料映射到低維度，同時保留原始資料的相對關係。其運算方法是計算兩筆資料之間的歐氏距離，然後透過矩陣運算獲得物體的位置並繪製出來。

以下是使用 MDS 繪製的人物相對位置圖。從圖中可以看出，雖然資料可以大致分爲左上、右上、中下三個群組，但各群組之間似乎仍然比較鬆散。爲了能夠更好的將各群區分出來，此處引入了一種相似度的算法。



圖十 MDS 分群效果

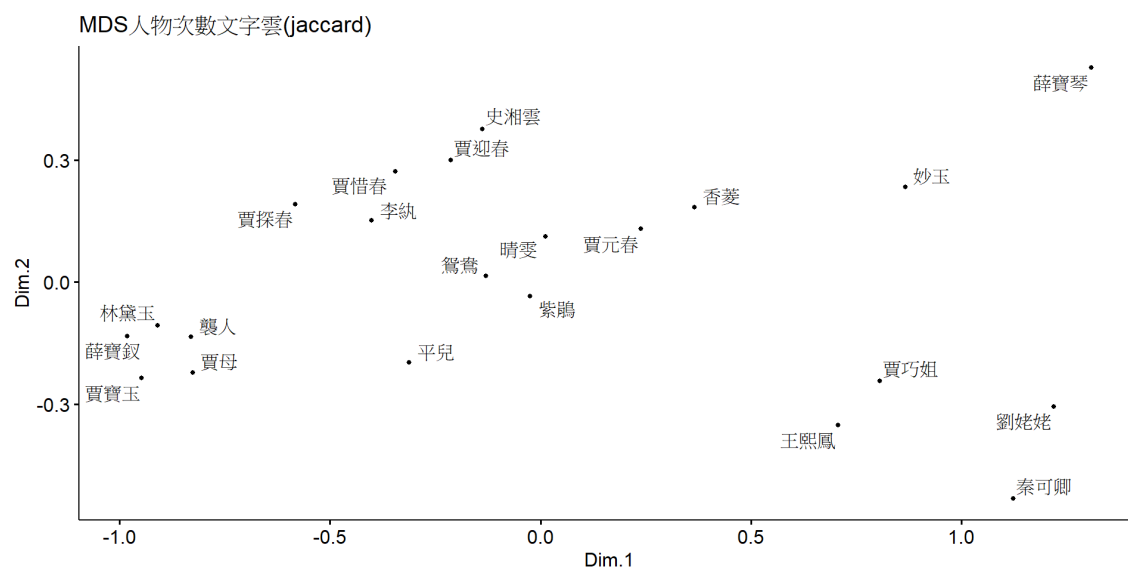
## 4.5 Jaccard

Jaccard 是一種衡量兩個物體相似程度的方法，其計算方式如式一，是為「兩個物體的聯集除以兩個物體的交集」。將 Jaccard 方法應用於上述 MDS 的資料中，便能顯示兩個人物在同一回中是否出現。

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

式一 Jaccard 計算方式<sup>[?] [?] [?]</sup>

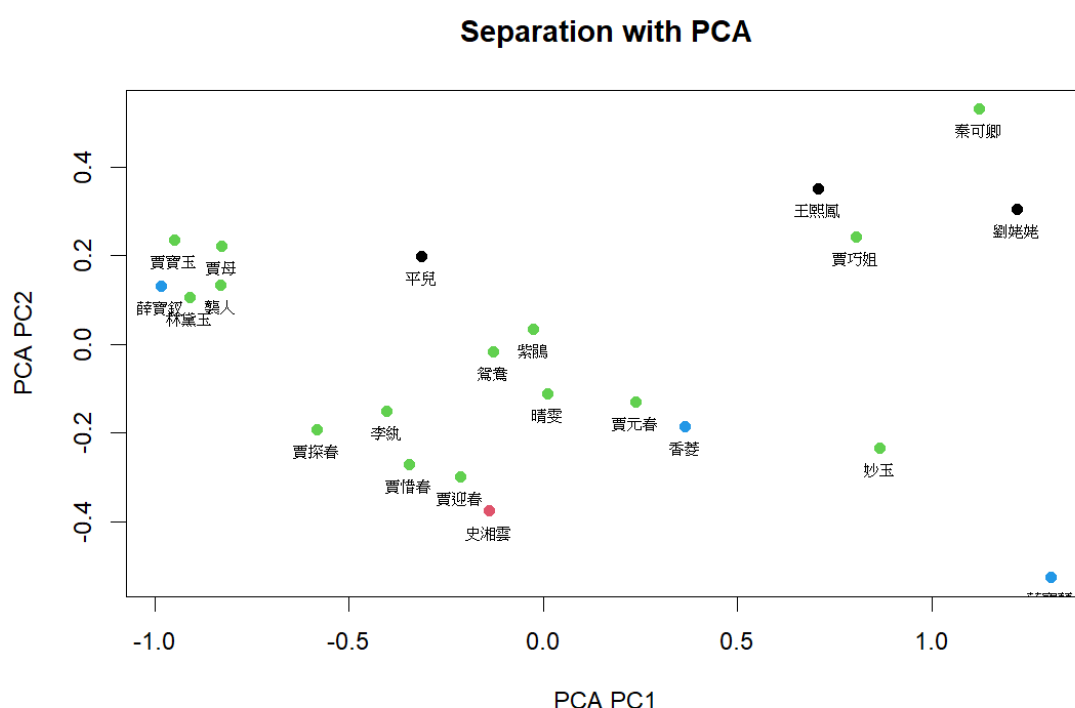
以下是添加了 Jaccard 後生成的 MDS 文字雲。從圖中可以清楚地看到群組中的點與點之間更加緊密，並且可以更好地觀察到人物之間的生活圈，以更好地判讀人物之間的關係。此外，藉由數據的相似度轉換，也可以有效將出場次數較少的人物分開，如 MDS 分群結果中妙玉看似和六到七位關係親密，但透過 Jaccard 轉換後，該人物即明顯與原先其他人物有一定距離，而小說故事狀況也是如此。



圖十一 MDS 經由 Jaccard 方法後的分群效果

## 4.6 PCA

同樣先將數據做 Jaccard 後使用 PCA 方法作圖後，發現其結果與 MDS 結果相似，且具有鏡射關係。依據圖中的遠近關係，概略分為左中右三大群，如賈母、賈寶玉等人為一群，賈元春、賈迎春等人為一群，剩餘將右側六位人物概括為一群，如賈巧姐、妙玉等。值得注意的是此次分群仍有些疑慮，如圖中右側的人物間仍具有一定距離，較難判斷人物之間是否有關係。因此我們考慮使用 t-SNE，嘗試使用另一種降維方法來檢視此方法對於人物分群是否更明顯。

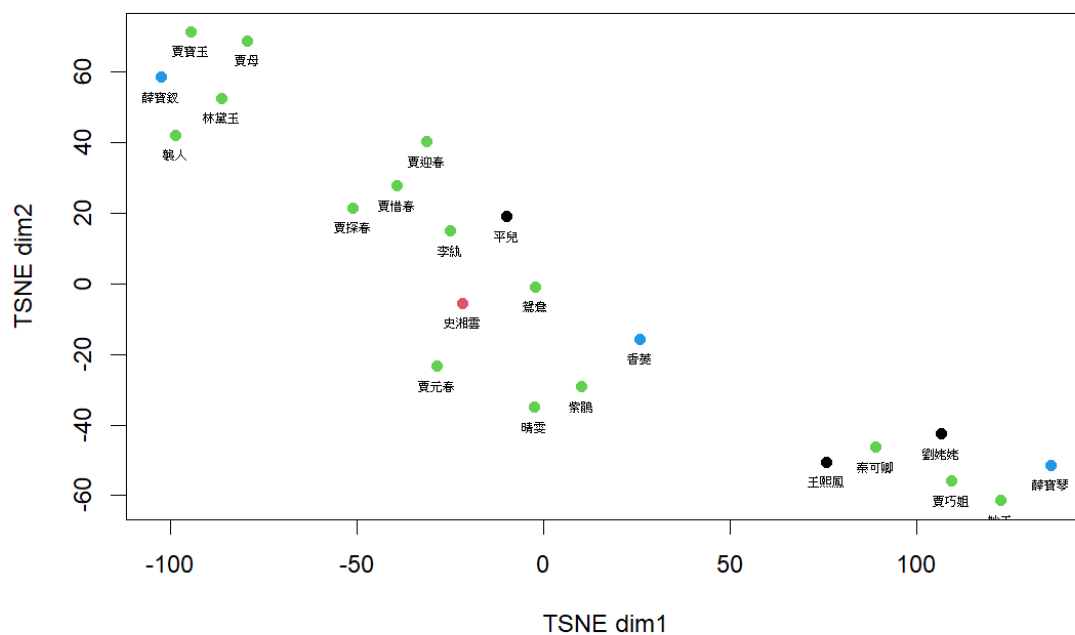


圖十二 PCA 分群效果

## 4.7 t-SNE

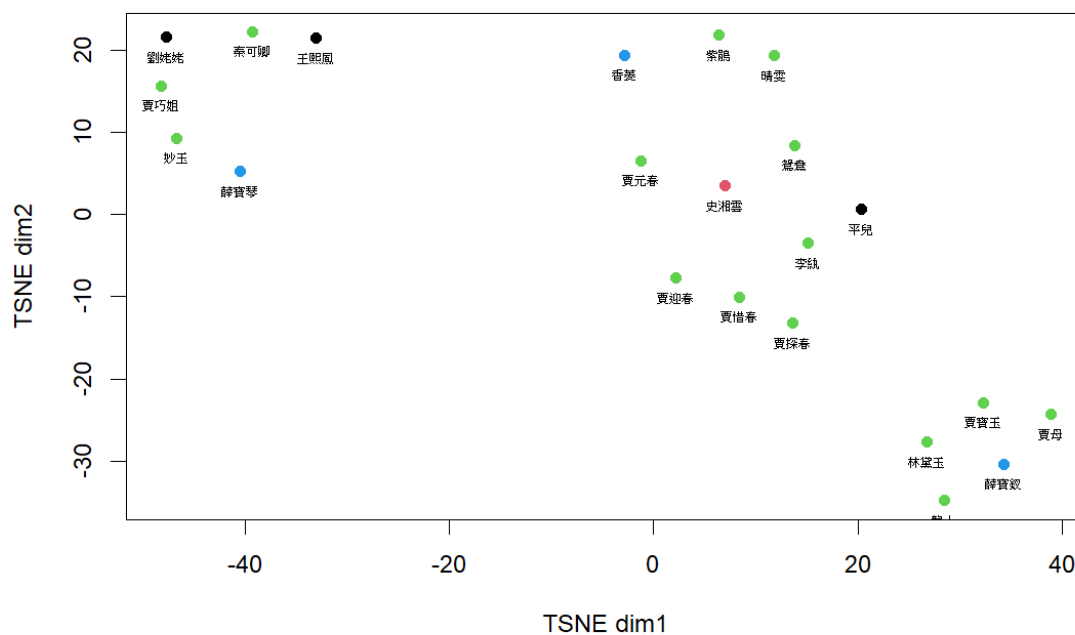
將數據藉由 t-SNE 方法降維，當困惑度 (perplexity) 設定於 5 後，發現相較於前兩個降維方法發生分群不明顯的狀況獲得改善，可以明確將圖中人物分為三大群，推測我們的數據較適合使用非線性降維方法，使得我們的模型得到更好解釋。值得注意的是，由於 t-SNE 降維後會產生不同的點分佈，故每次執行結果皆不相同，因此需要重複執行多遍。以下兩張圖是分群較為滿意的結果：

### Separation with the Rtsne package



圖十三 t-SNE 分群結果

### Separation with the Rtsne package



圖十四 t-SNE 分群結果

## 第五章 結論及建議

在決定撰寫以《紅樓夢》為主題的文字探勘時，考慮到該作品以白話文為主，且人物之間的關係複雜而密切，因此我們認為這是一個理想的文本進行探勘和分析的對象。

由於分析的數據屬於二進位制非對稱類型，我們將 Jaccard 方法應用於上述的數據中，以更精確地衡量關鍵字之間的相似性。在研究中，我們分別選擇線性降維方法(如 MDS 和 PCA)以及非線性降維方法(如 t-SNE)進行操作，將高維度資料映射至二維度的空間，使得圖形更加容易理解，同時也盡可能保留更多的訊息。結果顯示，t-SNE 應用在人物出現次數的數據時，是三者中最能夠有效區分不同群體的方法。

本研究的方法也可以應用於網路媒體和文本挖掘領域，與最近著名的 Chat-GPT 相似，都是致力於從文本數據中獲取有價值的訊息。另外，對於文學領域而言，本研究提供了一個相對客觀的工具，能夠自動化分析文本並提取重要資訊，例如人名、虛字等等。即使對文本內容一無所知或了解有限，研究結果也能讓人了解關鍵字之間的關聯性，不論讀者對該文本的內容是否具有深入了解，這個研究成果都能為他們提供一份有價值的資源，類似於一份「懶人包」。

研究過程中，我們還創建了人物總出現次數的數據，也就是《紅樓夢》中各人物在各回中的總出現次數。未來的相關研究可以借鑒我們的分析方式，比較和原模型差異之處，並建議探索其他降維方法，以尋找更合理的解釋模型。另外在蒐集資料時，也有找到其他研究者做過植物出現頻率的研究，因此未來可以針對其他主題進行關聯分析，找尋各項變數間彼此的關係。而在視覺化圖形結果方面也呈現出初步的懶人包效果，我們認為相較於單單的數據呈現，圖形化的結果對於觀眾而言會是較簡便的呈現方式。在未來我們希望能夠透過 Shiny 套件或其他更好的方法實現將資料結果以互動式視覺化網頁的方式呈現，讓所有的資料都能夠與使用者互動，從而增進互動上的便利性，更有利服務人群。

## 參考文獻

- [1] 曹雪芹、高鶚（約為乾隆初年）。紅樓夢。（程偉元編）。維基文庫。取自<https://zh.wikisource.org/zh-hant/>
- [2] 顏守玄（2021）。互動式資料視覺化之文字探勘以金庸三部小說為例 *Interactive Data Visualization of Text Mining to Jin Yong Three Novels*。臺灣博碩士論文知識加值系統。取自<https://hdl.handle.net/11296/9a2p75>
- [3] 喵喵（2020）。| 讀 | CN | 紅樓夢 | 人物關係表 | 四大家族 | 寶玉、黛玉、寶釵血緣關係 | 5 分鐘看懂紅樓夢人物關係 | 淺談紅樓夢 | 2023。NightelfMeowMeow。取自<https://nightelfmeowmeow.com/378/1105014-reading-cn-novel-dream-of-the-red-chamber-family-tree/>
- [4] 重明論（2021）。R 語言基於共現提取《雪中悍刀行》人物關係並畫網絡圖。知乎。取自<https://zhuanlan.zhihu.com/p/388637831>
- [5] Alan H. Lipkus. (1999, October). *A Proof of the Triangle Inequality for the Tanimoto Distance*. Journal of Mathematical Chemistry. 1999, 26 (1-3): 263—265 <https://link.springer.com/article/10.1023/A:1019154432472>
- [6] Flodel. (2012). *How to Create an Edge List from a Matrix in R?* Stack Overflow. <https://stackoverflow.com/questions/13204046/how-to-create-an-edge-list-from-a-matrix-in-r>
- [7] Finnstats. (2021, November). *How to Calculate Jaccard Similarity in R*. R-Bloggers. <https://www.r-bloggers.com/2021/11/how-to-calculate-jaccard-similarity-in-r/>
- [8] Michael Levandowsky, & David Winter. (1971). *Distance between Sets*. Nature. 1971, 234 (5): 34—35, <https://www.nature.com/articles/234034a0>
- [9] Sven kosub. (2016, December). *A Note on the Triangle Inequality for the Jaccard Distance*. <https://arxiv.org/pdf/1612.02696.pdf>
- [10] *Plotting PCA (Principal Component Analysis)*. [https://cran.r-project.org/web/packages/ggfortify/vignettes/plot\\_pca.html](https://cran.r-project.org/web/packages/ggfortify/vignettes/plot_pca.html)
- [11] *R: The R Project for Statistical Computing*. <https://www.r-project.org/>

- [12] *RStudio*. Posit. <https://www.rstudio.com>
- [13] *TSNE test*. <https://www.bioinformatics.babraham.ac.uk/tsne/>
- [14] *Wordlayout: Word Layout. In wordcloud: Word Clouds*. (2019). Rdrr.io. <https://rdrr.io/rforge/wordcloud/man/wordlayout.html>

# 附錄

## R 語言程式碼

本研究之 R 語言程式碼，使用 R 與 RStudio 製作。

<MINTED>

## 資料檔案

程式碼中需導入的資料檔案:

《紅樓夢》各章回（資料來源：維基文庫<sup>[?]</sup>）

紅樓夢\_第 001 回.txt

第一回 甄士隱夢幻識通靈 賈雨村風塵懷閨秀

【甲戌、庚、蒙批：此開卷第一回也。】

作者自：因曾歷過一番夢幻之後，故將真事隱去，而借通靈之說，撰此《石頭記》一書也。故曰「甄士隱夢幻識通靈」。但書中所記何事，又因何而撰是書哉？自又：今風塵碌碌，一事無成，忽念及當日所有之女子，一一細推了去，覺其行止見識，皆出於我之上。何我堂堂之鬚眉，曾不若彼裙釵哉！實愧則有餘，悔又無益之大無可奈何之日也！當此時，則自欲將已往所賴，上賴天恩，下承祖德，錦衣紈袴之時、飫甘饜美肥之日，背父母教育之恩，負師兄規訓之德，已至今日一事無成、半生潦倒之罪，編述一記，以告普天下人。我之罪固不能免，然閨閣中本自歷歷有人，萬不可因我之不肖，自護其短，則一併使其泯滅也。雖今日之茆椽蓬牖，瓦灶繩床，其風晨月夕，階柳庭花，亦未有傷於我之襟懷筆墨者。雖我未學，下筆無文，何為不用假語村言，敷演出一段故事來，以悅人之耳目哉。故曰「風塵懷閨秀」，乃是第一回題綱正義也。開卷即曰「風塵懷閨秀」，則知作者本意原為記述當日閨友閨情，並非怨世罵時之書矣。雖一時有涉於世態，然亦不得不敘者，但非其本旨耳，閱者切記之。

：

不知有何禍事，且聽下回分解。

表四 紅樓夢\_第 001 回.txt（節錄）

紅樓夢\_第 002 回.txt

紅樓夢\_第 003 回.txt

紅樓夢\_第 004 回.txt

：

紅樓夢\_第 120 回.txt

人物表.xlsx

name	type	group
賈寶玉	nr	賈家
林黛玉	nr	賈家
薛寶釵	nr	薛家
賈元春	nr	賈家
賈探春	nr	賈家
史湘雲	nr	史家
妙玉	nr	賈家
賈迎春	nr	賈家
賈惜春	nr	賈家
王熙鳳	nr	王家
賈巧姐	nr	賈家
李紈	nr	賈家
秦可卿	nr	賈家
賈母	nr	賈家
劉姥姥	nr	王家
香菱	nr	薛家
平兒	nr	王家
晴雯	nr	賈家
襲人	nr	賈家
紫鵲	nr	賈家
鴛鴦	nr	賈家
薛寶琴	nr	薛家

表五 人物表.xlsx

多稱謂人物.xlsx

賈寶玉	林黛玉	薛寶釵	賈元春	賈探春	史湘雲	妙玉	賈迎春	賈惜春	王熙鳳	賈巧姐	李紈	秦可卿
賈寶玉	林黛玉	薛寶釵	賈元春	賈探春	史湘雲	妙玉	賈迎春	賈惜春	王熙鳳	賈巧姐	李紈	秦可卿
寶玉	黛玉	寶釵	元春	探春	湘雲	妙玉	迎春	惜春	熙鳳	巧姐	李紈	可卿
此石	顰顰	衛蕪君	賈妃	蕉下客	枕霞舊友	檻外人	二木頭	藕榭	鳳辣子	妞妞	宮裁	蓉大奶奶
寶二爺	顰兒	寶姐姐	元妃	三姑娘			二姑娘	四姑娘	璉二奶奶	大姐姐	稻香老農	兼美
怡紅公子	林姑娘	寶丫頭	貴妃									秦氏
絳洞花王	林丫頭		大姑娘									
	瀟湘妃子											

表六 多稱謂人物.xlsx