



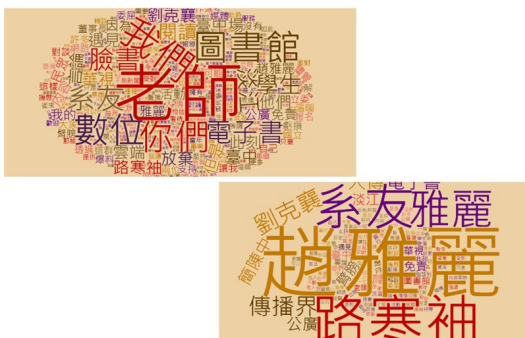
「教師跨領域研究社群」-文字探勘與資料庫建置

時間：2022-11-11

地點：傳播館 Q403

唐大崙分享自己過去耗費數年，從網路上爬取新聞文章，所建立的 300 多萬則新聞資料庫中，辛苦建立出來的、近 60 萬個中文詞的統計特徵資料庫，包括 TF、IDF、分布標準差等資料，依據此資料庫便能輕鬆判斷一則新聞的關鍵詞究竟是哪一些詞彙，也可以依據關鍵詞與詞彙所在相對位置的權重調整，計算兩篇文章相似度。這種相似度計算結果，與人的主觀判斷相當接近，表示這個演算法可以被接受用來追蹤語意相似的文章。不過在語意追蹤計算之前，需要預先建立正確斷詞的過度檔案，這是比較複雜的過程。但是，這樣的資料庫是否適用於法律契約的文章？仍是個必須真實嘗試才能確定的實徵性問題。因此，唐大崙與蔡明修老師約定下次再進一步拿真實契約文章來嘗試，也討論到先合作開設一個免費為人計算中文文章相似性的服務網站，等技術更精進，再進一步客製化為建築契約文件做語意追蹤的服務，這將可開啟另一種跨領域的合作範例。

詞頻文字雲 與 詞關鍵度文字雲 的差別



唐大崙以同一篇新聞做詞頻分析 對比於 關鍵詞分析所做的文字雲顯示，有很大差異。
唐大崙老師先分享文字探勘部分的經驗與做法。



唐大崙說明多篇文章之間的相似度資料可用以進行像這類分群的語意追蹤網絡圖。
唐大崙談到要計算出關鍵詞需要像這樣的中文字詞分布的統計資料庫。